
Cognitive Goals as Reward-Producing Programs

Guy Davidson ^{α *}, Graham Todd ^{β *}, Julian Togelius ^{β} , Todd Gureckis ^{γ} , Brenden M. Lake ^{α, γ}

^{α} Center for Data Science, ^{β} Game Innovation Lab, ^{γ} Department of Psychology
New York University

guy.davidson@nyu.edu, gdrTodd@nyu.edu

Abstract

People have a remarkable capacity to generate their own goals, beginning with child’s play and continuing into adulthood. Despite considerable empirical and computational work on goals and goal-oriented behavior, models are still far from capturing the richness of everyday human goals. Here we bridge this gap by collecting a dataset of human-generated playful goals, modeling them as reward-producing programs, and generating novel human-like goals through program synthesis. Reward-producing programs capture the rich semantics of goals through symbolic operations that compose and add temporal constraints, and allow for program execution on behavioral traces to evaluate progress. To build a generative model of goals, we learned a fitness function over the infinite set of possible goal programs, and sample novel goals with a genetic algorithm. Human evaluators found the model’s better samples indistinguishable from human-created games. We also discovered that our model’s internal fitness scores predict games that are evaluated as more fun to play and more human-like.

Understanding how humans create, represent, and reason about goals is crucial to understanding human behavior. Goals are pervasive throughout psychology [19, 3, 21], having been studied from perspectives such as motivation [33, 20, 6], personality and social psychology [23, 56], and learning and decision-making [50, 49]. But what is a goal? Elliot & Fryer present one workable, albeit simplified definition: “A goal is a cognitive representation of a future object that the organism is committed to approach or avoid” (see also [21, 49]). Reinforcement learning offers another formulation, operationalizing goals as maximizing cumulative reward over a series of steps [63]. Typical goals in reinforcement learning tasks include reaching a target location, winning in a video or board game [48], or placing an object in a specified position (e.g. Figure 1a), such that success can be characterized by reaching a target state.

In contrast, people routinely create novel, idiosyncratic goals with a richness that goes beyond these common modeling settings. Chu et al. [14] report the example of Gareth Wild, who set an unusual goal for himself to park in every spot in a particular grocery store’s parking lot (Figure 1b). Children routinely devise fun and compelling goals without external guidance, such as creating a “truck carrier truck” (Figure 1c) or stacking as many blocks as possible in a single tower (Figure 1d). Beyond being fun, these playful goals play a crucial role in learning to structure and solve arbitrary problems [13, 43, 1]. Indeed, the capacity to autonomously set and achieve goals is arguably a core component of human intelligence [14, 53].

We propose a framework for modeling human goal generation as synthesizing reward-producing programs (Figure 1, bottom row). There are several advantages to representing goals as symbolic programs, mapping an agent’s behavior to a reward score indicating the degree of success. First, a structured language facilitates the compositional reuse of motifs across disparate goals. This, in turn, makes capturing the wide range of human creativity in goal creation substantially more tractable: in Figure 1e, we illustrate a simple ball-throwing game (in black) and four distinct variants (in red, blue, pink, and brown) composed in part from shared components: balls being thrown (highlighted in

yellow), the thrown ball hitting something (orange), and the thrown ball landing somewhere (green). Second, our choice of representation makes goal semantics explicit. The particular grammatical elements of our representation each fulfill particular roles, such as *predicates* (i.e. specific and evaluable relations between objects, colored orange in the programs in Figure 1) and *temporal modals* (i.e. relationships in time between goal components, such as ‘until’ and ‘then’ in Figure 1). Finally, goals-as-programs are executable; that is, they can be computationally interpreted to detect when a goal is entirely or partially achieved (Figure 1e, each program would be interpreted and provide a score only when the matching throw trajectory is achieved).

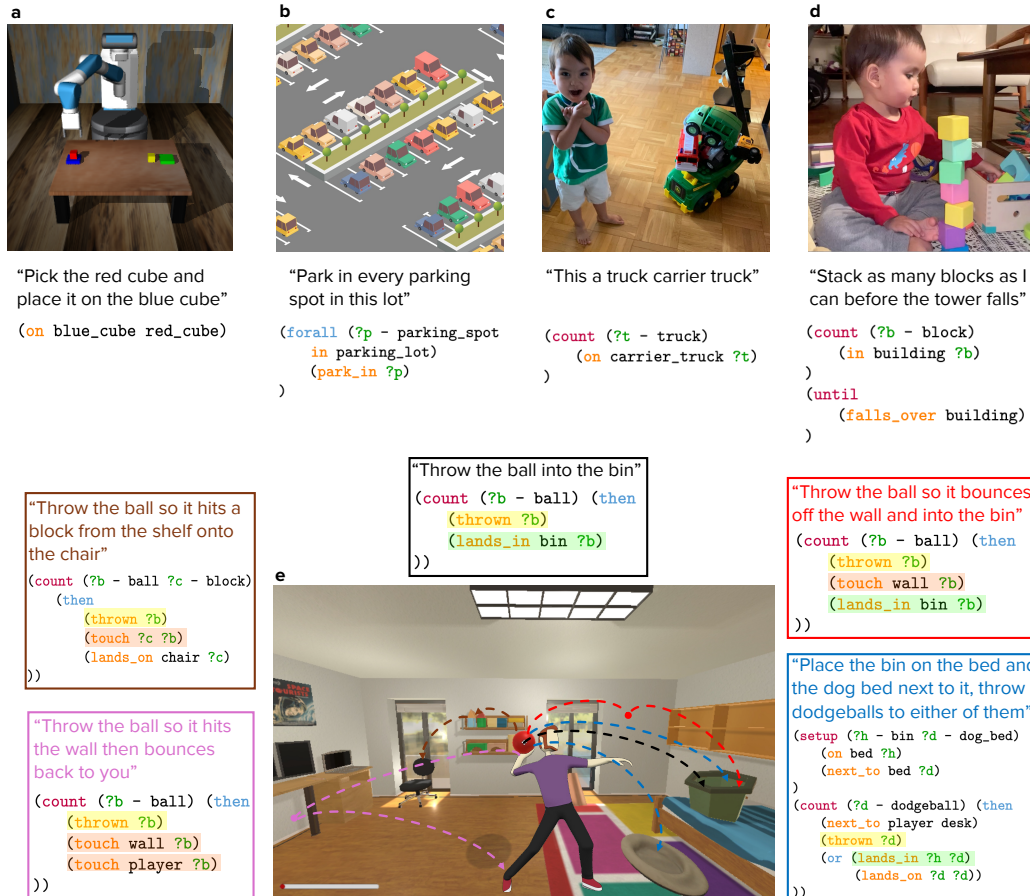


Figure 1: Generating Human-Like Goals by Synthesizing Reward-Producing Programs. Panels a-d each represent a different goal, presented in natural language and mapped to pseudo-code in a program-like representation. Panel e represents several plausible goals in our experiment environment, some of which were created by participants in our experiment. Each goal is represented by a throw trajectory (dashed line in the illustration) matching a description of the goal (whose text is the same color as the line). We highlight shared compositional components between programs in yellow, orange, and green. Our program representations are reward-producing, that is, run on sequences of agent interactions with an environment (state-action pairs) and emit a score with respect to the specified goal. Our pseudo-code and domain-specific language both use a LISP-like syntax, with operator invocations in prefix notation, as the first argument inside the parentheses. Some of these goals were created by participants in our experiment; see Figure SI-1 for representations of the blue and pink programs in our domain-specific language.

In this article, we demonstrate that programs can capture real human-created goals in a specific domain and build a model capable of generating new programs representing human-like goals. Although goals play an important role in psychological theory, there are few, if any, empirical paradigms for eliciting goals from study participants. To that end, we collected a dataset of playful goals (games) created in our experimental environment. We then constructed a concrete program representation for these goals by creating a custom domain-specific language. Next, we analyzed our dataset of goal programs and found that consistent with past work on creativity [69], human games in our domain

are interesting and diverse while also sharing many common features. Participants also leveraged common sense about the physical world in how they used various objects in their goals.

With this (small) dataset of goal programs and the grammar for the programming language, we next developed a computational model to generate new goals in this representation. We view this as an important test of the framework we propose as generating new exemplars is a core capacity of representing and understanding concepts [69]. This task offers a particularly meaningful challenge as generating novel, human-like goals is a difficult, open-ended task that requires capturing various human cognitive capacities (such as common sense and creativity). Our model learns a fitness measure predicting human likeness from the program representations and then samples new programs that score highly on this metric. We found that our model succeeds in generating novel games distinct from examples in the training dataset. We then conducted a human evaluation of model-generated games with a new set of participants and asked them to rate several game characteristics including how “human-like” they were. We found that model games from sections of program space closer to participant-created games were judged indistinguishably from the real games, but some model samples further away were not rated as highly. Analyses revealed that our learned fitness function predicts several of the human judgment questions, including how fun and human-like games are found, giving an indication that our goal representations capture some aspects of how people creatively construct new goals.

Behavioral results

We created an experimental setting that would allow us to capture the rich, playful, and creative nature of how children (and adults) create goals in everyday scenarios. We used AI2-THOR [37] (an embodied, 3D environment simulation) to set up a room resembling a child’s bedroom, filled with toys and other common objects. In our task, we asked subjects to propose a single-player game to be played in the room. This simple, open-ended design allowed participants to imagine and propose a wide range of playful goals, with the aim of game generation helping to make the resulting goals more concrete. We collected a dataset of 98 games, described by participants in natural language. In addition, we recorded full state-action traces of each participant’s interactions with the environment, which we leveraged in later experiments (see Dataset collection methods for additional details).

We then translated each game from natural language to a program-like representation in a domain-specific language (DSL). Our approach is inspired by instantiations of Fodor’s Language of Thought using programs [25, 27, 61, 70], and consistent with the framework we propose in the article, we use programs to explicitly represent meaning (semantics) separately from linguistic minutia. Our DSL is designed to capture the key semantic components of games in our dataset and is initially derived from the Planning Domain Definition Language (PDDL [26]), which offers a basic representation for specifying goals (i.e. end states of plans) and preferences (i.e. other costs to optimize while planning). Each program in the DSL contains two mandatory sections: gameplay **preferences** describing how a game is played, and **scoring** rules specifying how to determine a player’s score based on the satisfaction of the game’s preferences. Game programs may also contain optional **setup** instructions and **terminal** conditions (see Supplementary information G for the full DSL). We manually translated the 98 games in our dataset to this DSL.

Our choice to represent games in a form with explicit semantics allows us to quantitatively analyze their structure and fundamental components. We found that people recruit an intuitive physical common sense when creating games (Figure ED-2, and see Game dataset analyses methods for details). If an object is thrown, it’s likely a ball, and if an object is stacked, it’s likely a block — and while a few participants specified games involving throwing blocks, none attempted to stack balls. Similarly, participants did not specify throwing cumbersome objects (such as the laptop or chair), and a participant who specified throwing a large ‘beach ball’ clarified that it should land *on* the bin (as the ball does not fit within the bin). We also observed evidence of both compositionality (common structure reuse) and creativity (preponderance of unique structures), summarized in Figure 2 (see Game dataset analyses methods for details). Counting occurrences of grammatical structures while abstracting over the identities of individual objects (i.e. treating (`once (agent_holds block)`) and (`once (agent_holds ball)`) the same), we find the five most common structures cover almost half of the total observations, offering strong evidence of reuse through compositionality. At the other end of the distribution, we also observe a long tail emblematic of creativity, as one-half of the unique structures we count appear exactly once. Despite not explicitly being encouraged

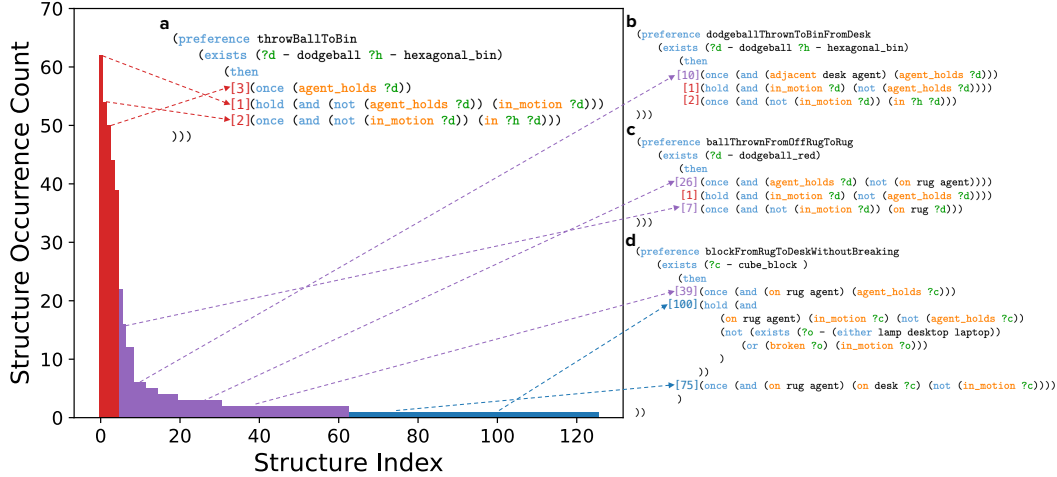


Figure 2: Participants showcase compositionality and creativity. We analyze the occurrence of various abstract structures in our programs (see Game dataset analyses methods for details). **Red:** The five most common structures cover almost half (47.5%) of total occurrences, showing a strong degree of compositional reuse. The three most common structures combine to create a simple ball-to-bin throwing preference (**a**, with structure indices in square brackets). **Purple:** Many other structures are reused fewer times, covering a large portion of the remaining occurrences (another 40.5%). Using some of these rarer structures allows for creating more complex throwing elements, such as limiting where the player throws the ball from (**b**) or changing the throwing target (**c**). **Blue:** Exactly half of the structures (63 / 126) appear only once — this long tail of expressions offers evidence of creativity. The last throwing preference (**d**), specifying throwing a block from the rug onto the desk without moving off the rug or breaking any of the objects on the desk, uses two unique structures.

to generate novel or creative games, many participants proposed entirely unique gameplay ideas, encouraging us in our endeavor to model the richness and creativity of human goals.

Modelling Results

To test the representational our framework, we develop a computational model that synthesizes human-like goals as programs. Our approach (illustrated in Figure 3 draws primarily on contrastive learning and evolutionary techniques, and our motivations are two-fold. As early as childhood, humans create novel goals without anything resembling extensive data-driven training. We therefore similarly use a realistic (i.e. not internet-scale) number of training examples. We also attempt to capture the key attributes of human goal generation by explicitly integrating some of the cognitive capacities (i.e. common sense and compositionality) that our previous results suggest people recruit in creating goals.

The first component of our model is a fitness function $f(g) = \theta \cdot \phi(g)$ that maps $f : \mathcal{G} \rightarrow \mathbb{R}$ from a game $g \in \mathcal{G}$ to a real-valued score that encodes its “human-likeness.” We transform each game into an 89-dimensional vector of features that capture properties relating to structure (e.g., the size and depth of its syntax tree) or logic (e.g., whether any expressions are redundant). Other features proxy cognitive capacities, such as physical common sense (estimating predicate feasibility from play data) or compositionality (n-gram model over syntax elements). We leverage our programmatic representation of goals in order to automate this feature extraction process (see Fitness function methods for details). We then learn weights θ for the features in a *contrastive* fashion [11, 39] by optimizing for the difference in scores between our set of human-generated games and a substantially larger set of corrupted (i.e. lower quality) games obtained through random tree-regrowth [27] on our dataset (see Figure ED-3a and details in Fitness function methods). This learned fitness function then guides an evolutionary search procedure in order to generate novel games. Broadly inspired by work in genetic programming, we use a *quality-diversity* algorithm [57, 10] called MAP-Elites [51] to generate a set of samples that widely cover the space of programs in addition to optimizing the fitness function (see Figure ED-3b). The details of our implementation, including the particular criteria used for maintaining sample diversity, are available in MAP-Elites methods.

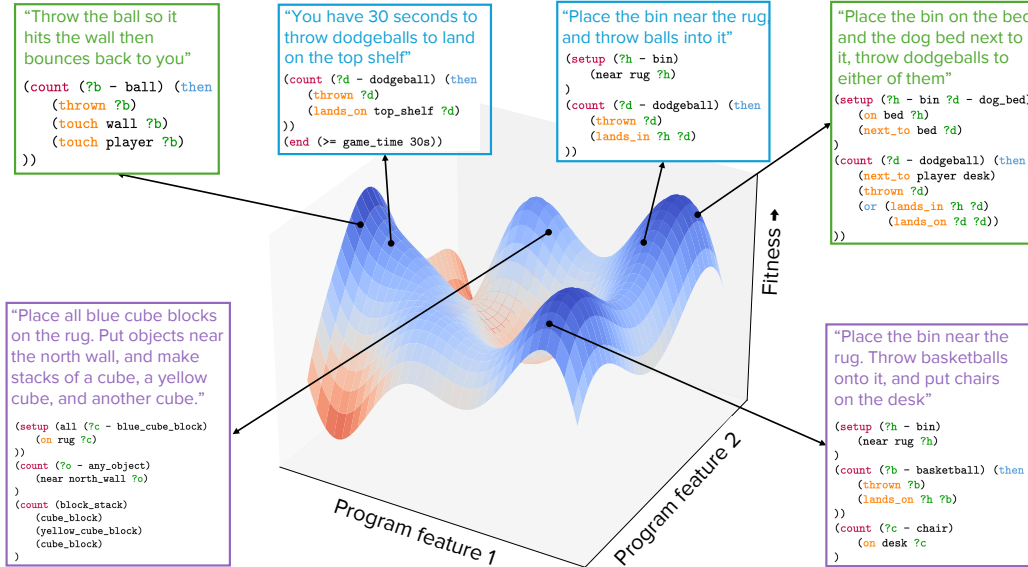


Figure 3: Computational model overview . Our model operates on programs in some high-dimensional space (visualized in two dimensions). We learn a fitness metric (Z-axis) capturing desirable aspects of programs, using a dataset of human-created goals (highlighted in green). Our model then generates diverse new samples maximizing the fitness measure, some “matching” participant-created goal programs on diversity criteria (in blue) and other “unmatched” novel goals (in purple). All goals in this figure were created by participants in our experiment or by our model; see Figure SI-1 for their representations in our domain-specific language.

Generated games

Quantitatively, our approach is successful: Figure ED-4 shows that the model quickly produces games with fitness scores in the range of human samples and does so across many of the “niches” defined by our behavioral characteristics. Considering only games that pass a plausibility characteristic we include in our MAP-Elites search, 1889 games (94.45%) exceed the fitness score of the least fit real game, and exactly 1000 games exceed the fitness of the median human game. To the extent that our fitness function captures human likeness, our model produces human-like games. Qualitatively, we find a variety of outputs that range from variants of simple games in our reference dataset to games in entirely new regions of program space. In Figure 4, we present examples of games generated by our model that “match” a game in our human dataset by occupying the same niche defined by the MAP-Elites algorithm. Overall, the general structure of matched games is similar (e.g. both the left-most human game and its corresponding matched generated game involve two throwing-based objectives). However, minor elements in generated games tend to be less readily explainable (e.g. the scoring condition in the middle generated game, which arbitrarily multiplies the number of satisfactions by 0.4). Our model also produces games that occupy niches without corresponding human games (Figure 5). While these “unmatched” games also use predicates in reasonable and novel ways, they tend to lack coherence at a more global level. For instance, the individual preferences defined by the center game in Figure 5 (picking up triangular blocks and throwing balls) make sense individually but do not readily combine to make an easily explainable game.

Human evaluations

To systematically and extrinsically evaluate our model, we performed human evaluations using a second set of human participants ($n = 100$; see Figure ED-5 for the evaluation interface and Human evaluation methods for details). Evaluated games belonged to one of three different categories: real games created by participants in our first experiment, matched games (one corresponding to and occupying the same niche as each real game), and novel unmatched games (games in Figure 4 and Figure 5 were included; see Human evaluation methods for details). Participants evaluated three games in each category above (without knowing their categories) in a randomized order and answered seven multiple-choice questions for each, each on a different measure, such as human likeness, fun,

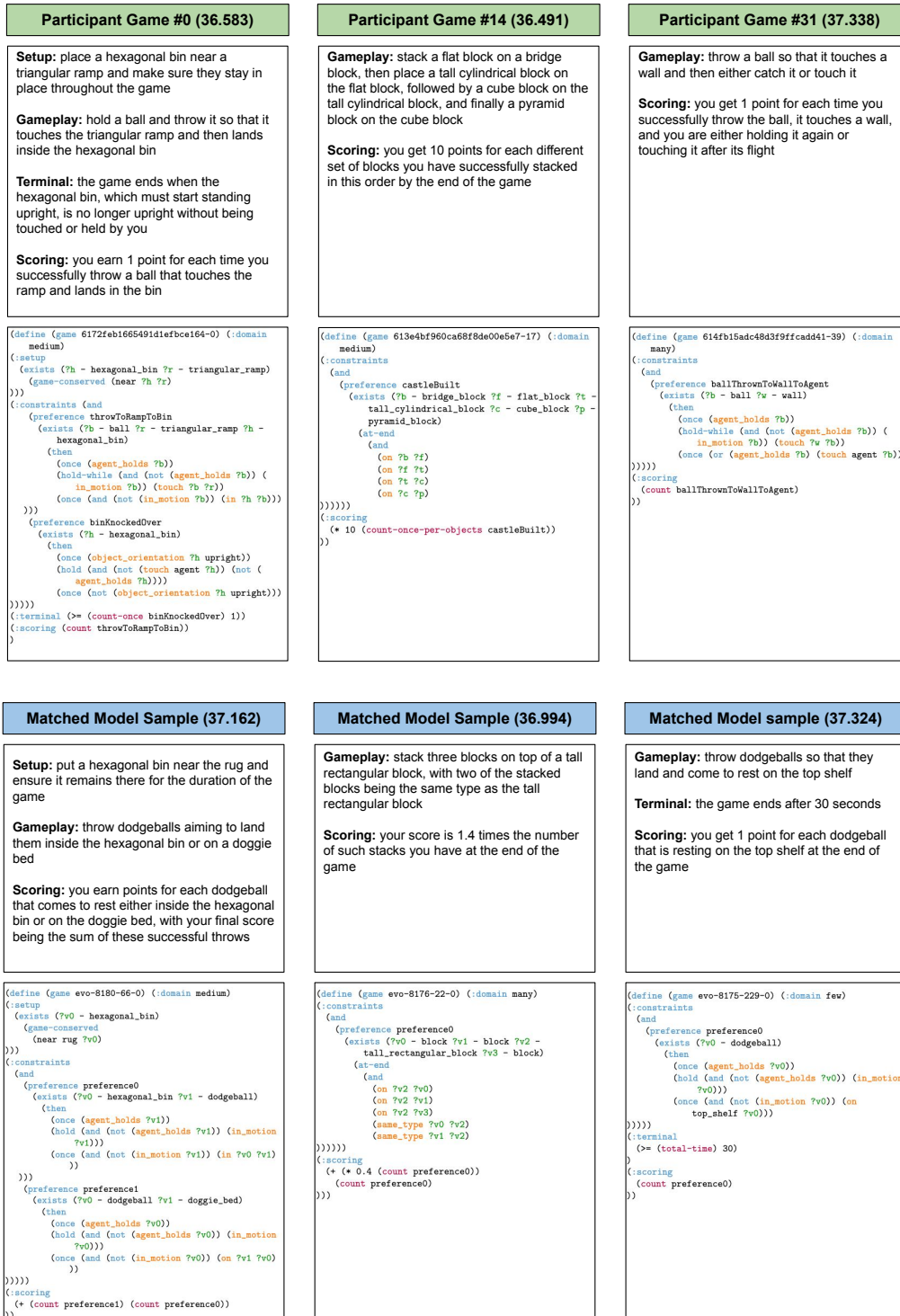


Figure 4: Our goal generation model generates simple human-like games. Each pair of games in a column has the same set of MAP-Elites behavioral characteristics. Parentheses: the fitness score assigned by the model to each game. Natural language descriptions are generated through automated back-translation from programs (see Supplementary information D for details)

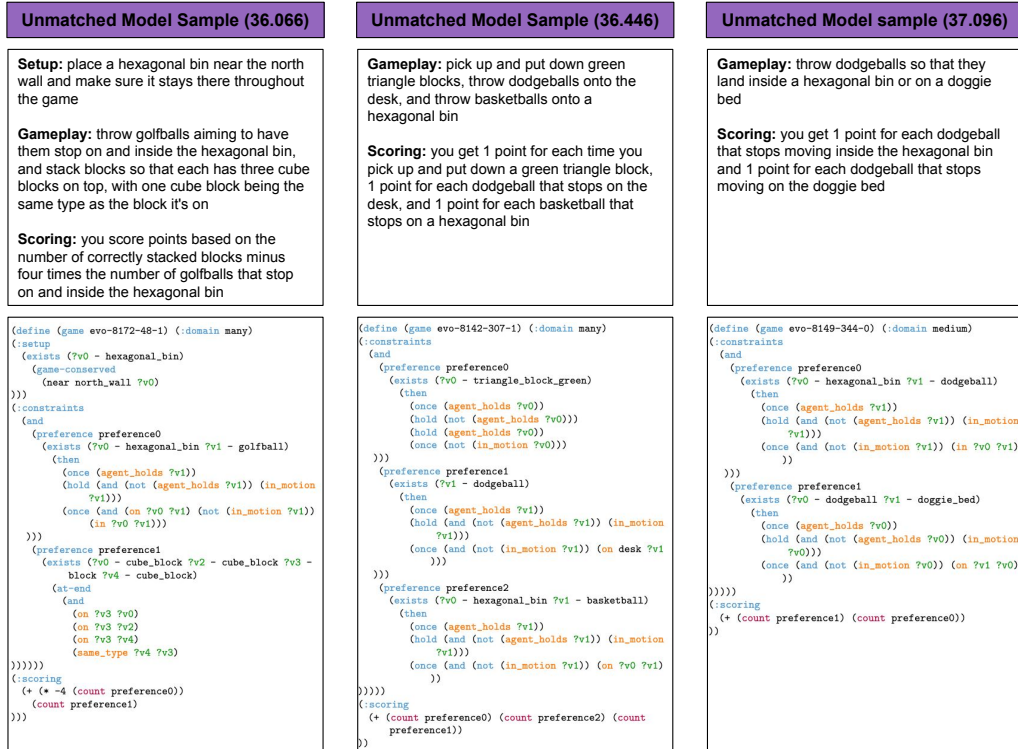


Figure 5: MAP-Elites generates interesting, novel goals. Each of the three games below has high fitness and fills a cell in the MAP-Elites archive with no corresponding human game in our dataset. Parentheses: the fitness score assigned by the model to each game.

and creativity. Our final evaluation dataset includes 892 participant-game evaluations, each with a response to all seven measures.

Table 1: Human evaluation result summary

| Measure | Mean score by category | | | Significance of difference | | |
|-------------------------|------------------------|-------------|---------------|----------------------------|--------|--------|
| | Real (R) | Matched (M) | Unmatched (U) | R vs M | R vs U | M vs U |
| <i>Understandable</i> ↑ | 3.943 | 3.923 | 3.331 | - | *** | *** |
| <i>Fun to play</i> ↑ | 2.522 | 2.430 | 2.068 | - | *** | *** |
| <i>Fun to watch</i> ↑ | 2.385 | 2.313 | 2.024 | - | *** | ** |
| <i>Helpful</i> † ↑ | 2.997 | 2.987 | 2.840 | - | - | - |
| <i>Difficult</i> ↓ | 2.582 | 2.660 | 2.676 | - | - | - |
| <i>Creative</i> ↑ | 2.318 | 2.213 | 2.143 | - | * | - |
| <i>Human-like</i> ↑ | 2.813 | 2.670 | 2.119 | - | *** | *** |

Evaluators don’t distinguish between participant-created real and matched model games, but do distinguish games from both. Participants responded to seven Likert questions on a 5-point scale, one for each attribute in the first column (see Human evaluation methods). We report the nonparametric Mann-Whitney U test [46] for differences in outcomes. See Supplementary Table SI-2 for test statistics and P-values. *: $P < 0.05$, **: $P < 0.01$, ***: $P < 0.001$.

†: The full measure description is “Helpful for interacting with the simulated environment.”

In most measures, higher scores are better, indicated by the ↑, other than *Difficult* ↓, in which 3 means “appropriately difficult”, and scores below and above indicate too easy and too hard respectively.

Table 1 summarizes the quantitative responses from our human evaluations. We begin with a simple statistical comparison of the ratings of the games in the different categories using the nonparametric Mann-Whitney U -test [46] (and see Human evaluation methods for additional details). Participants respond similarly to the real and matched games, with no statistically significant differences in

average response scores across all seven attributes. On the other hand, the unmatched games differ on a number of attributes. Compared to both real and matched games, participants judge them to be less easily understood, less fun to play and watch, and less human-like.

To further analyze these differences and the extent to which they are mediated by our fitness measure, we performed a mixed-effects regression analysis whose results are summarized in Table ED-1. We fit independent models using each of the seven attributes we have our human evaluators judge as the dependent variables. We include fixed effects for the fitness score and indicator variables for membership in the real and matched groups (treating the unmatched group as a baseline), and random effects for the participants and individual games (see Human evaluation methods and Supplementary Table SI-3 for full details). We find that our fitness function captures many of the evaluated attributes: higher fitness predicts higher ratings of understandability, fun to play, and human likeness ($\beta_{\text{fit}} > 0$); conversely, higher fitness also predicts lower ratings of helpfulness, difficulty, and creativity ($\beta_{\text{fit}} < 0$). We find the positive predictability promising: that our fitness function, learned to maximize human likeness in a symbolic program space, also captures notions of understandability and fun. Conversely, we view the negative relations as evidence of some degree of mode-seeking: our fitness measure likely assigns the highest scores to the games most representative of the dataset at large. These modal games are plausibly neither particularly creative nor difficult, which means that participants might find also them less helpful. However, even accounting for any mediating effect of fitness scores, differences in attribute ratings between game groups persist (see Supplementary information E.2 for details). Given that these tests indicate that participants respond similarly to matched and real games, one potential explanation is that model-generated samples simply replicate the corresponding human-generated ones. We examined this question and found that these games are substantially functionally different (see summary in Figure ED-6, details in Supplementary information E.4 and methodological details in Sample similarity comparison methods).

We also performed ablations of key model components corresponding to the cognitive capacities we found our participants recruit. Our model operationalizes *physical common sense* through two fitness features that estimate predicate expression feasibility using our database of participant environment interactions. If we lesion these two fitness features, the resultant model generates goals with lower fitness and gameplay elements that are less likely to appear in our database of participant interactions (see Supplementary information F.1 for details). Our program representations are inherently highly compositional, making ablating the role of *compositionality* in our model challenging. We do so via lesioning the crossover mutation operator which explicitly leverages compositionality, and find that removing it drastically reduces sample fitness (see Supplementary information F.2 for details) Humans intuitively create *coherent* games, generating playful goals in which rules and components refer back to one another. Our model, which improves samples gradually using evolutionary operators, has more difficulty in maintaining global coherence. Like with common sense, our model operationalizes coherence through a few fitness features. Lesioning these features leads to model samples that are similarly worse, as measured using our database of participant interactions (see Supplementary information F.3 for details). We note that these features mostly capture lower-level coherence between fragments, but improving higher-level coherence remains a challenge for future work.

Discussion

Goals are a critical aspect of human cognition and, in fact, the starting place for many models of human behavior. However, the representation of goals is often impoverished. In this article, we proposed a new framework for understanding human goals as reward-producing programs, and thus understanding goal generation as program generation. To evaluate this framework, we developed an interactive experiment in which participants created playful goals, operationalized as games to be played in a virtual environment. By analyzing the program-based translation of these games, we highlighted several cognitive capacities recruited by our participants, such as physical common sense and compositionality. These capacities, in turn, helped inform our modeling efforts. We then built a computational model that learns from a small dataset of games and generates goals that are both novel and human-like according to human evaluators.

This work unites various strands of research in cognitive science, artificial intelligence, and game design. First, we build on a substantial literature studying the psychology of goals [19, 3, 21, 49, 14].

Our work also relates to goal-conditioned agents [45], and specifically *autotelic* agents capable of setting their own goals [15]. Typical goals in such systems are derived from agent or object positions [24, 8, 64], natural language descriptions [32, 18, 16, 68], or limited temporally-aware mechanisms [44, 35, 22, 42]. However, to date, these approaches lack the variety and richness of human-created goals. Finally, we draw inspiration from the automatic game design literature, such as synthesizing board game variants [55, 30, 7] or simple video games [66, 62, 71, 36]. Unlike our approach, these efforts often optimize program synthesis for some heuristic notion of “fun” [7, 66] rather than explicitly modeling human-like game generation.

In instantiating our model, we necessarily make a number of specific implementational and algorithmic decisions. Our framework is committed to the representation of goals as *reward-producing programs*: computationally executable mappings from agent behavior to indications of progress towards a goal. We find it crucial that these programs capture the rich, temporally extended nature of goals people create, and that they facilitate the flexible and compositional creation people seem to engage in [38, 69]. Our model operates over these programs and assumes two primary components: a measure of human likeness and a method to sample novel programs. However, supporting the specifics of this program language will require further work. In addition, this project currently manually translates from natural language descriptions to the proposed mental language of goal programs, a step inducing potential arbitrariness in how to perform this mapping. Future work will leverage recent advances in text-to-code language models to automate this step (see [70] for a discussion on using language to construct meaning through programs). Another key detail is our approach to sample diversity, which arises from our choice of “behavioral characteristics” that define the axes along which the MAP-Elites algorithm maintains diversity. In this work, we select behavioral characteristics based on notable gameplay components observed in our human dataset; future work could explore other techniques for maintaining diversity, including the automated selection of behavioral characteristics [17, 29].

We conclude by highlighting some promising avenues for future work. Our model, though powerful and productive in this domain, is inherently coupled to the environment and dataset we collected — particularly given the engineering effort to instantiate various types of knowledge. This approach has some distinct advantages: we can isolate various cognitive capacities, interpret their contribution to our fitness measure (Supplementary information B.1), and ablate their roles (Supplementary information F). Simultaneously, some of the challenges our model faces (such as coherence between program components) might be alleviated by incorporating natural language or by leveraging the capabilities of large language models to write code and adapt to in-context instructions. Future work could productively investigate approaches that require less tailoring to specific environments, which could then be leveraged to improve the abilities of reinforcement learning agents to follow richer goals across different domains. Our current model generates goals to maximize a single fitness measure, but many objectives have been suggested to guide human learning, such as curiosity [4], information gain [60], or learning progress [65]. It would be fruitful to instantiate goal generators to maximize these various objectives and compare the resulting goals and the behaviors that arise in artificial agents that pursue them.

Finally, we view our work in this article as a proof of concept for the broader framework we propose for understanding goals as programs. The model we instantiate operates at Marr’s [47] *computational* level; it tackles the *what* of goal generation without directly inspecting *how* people do so. We are excited for future work to seek algorithmic-level signatures and constraints on how children and adults flexibly generate rich goals and leverage them to build models that operate closer to Marr’s *algorithmic* level. Our formalization of goals as reward-producing programs makes explicit both the semantics of goals and their function to evaluate goal satisfaction. We believe such goals can help build richer *autotelic* [15] that pursue more diverse creative and human-like goals, and are keen to pursue this direction in future work. While many important questions remain open, we see palpable promise in improving the understanding of human goals through the framework we propose, and we hope this article provides evidence of the potential and utility of our approach.

Methods

Dataset collection methods

Experimental design: After a consent form and instructions quiz, participants completed a tutorial designed to familiarize them with the controls for our environment. After successfully completing the tutorial, they were placed in one of three variations of the main experiment room, with the same structure but different amounts of available toys and objects. Participants were then free to explore this new room until they had a game ready, and could freely reset it to its initial state in the meantime. Participants were asked to create games with the following restrictions: single-player, require no additional space or objects that they do not see in the room, and include a scoring system. While the latter constraint seems limiting, we note that any arbitrary goal can be scored by rewarding the achievement of the goal.

Dataset collection: Participants then reported their game in natural language in three text boxes, one of which was optional (see Figure ED-1). The optional first one allowed specifying whether there was any setup or preparation required to get the room from its default initial state to one that would allow playing the game (e.g., placing the bin on the bed). The second text box allowed participants to describe the game’s gameplay, and the third offered space to describe any scoring rules. To encourage participants to imagine playing their game, they were also asked to report their perceived difficulty level and how many points they thought they might score if they played it. Participants then had a chance to play their game and revise it should they want to; if participants opted to revise their games, we analyzed the revised ones. We contacted 192 participants via Prolific [54] of whom 114 finished the experiment and another 12 were paid due to technical difficulties. Participants were paid a base rate of \$10 and received a \$2 bonus if their game satisfied the required constraints. Successful participants took 44.4 minutes on average, with a standard deviation of 23.3 minutes. We then excluded 8 games that did not satisfy the constraints we posed on participants, 6 duplicates (including some due to technical difficulties from participants who restarted the experiment), and 6 other games that were unclear or under-specified. After accounting for two other games we opted to avoid modeling due to their complexity, we arrived at our final dataset of 98 games.

Interaction traces: In addition to the game descriptions in natural language, we record traces of participants’ interactions with the environment. We record state-action traces to allow us to replay and examine how participants interact with our environment. We record separate traces for each different segment of the experiment (before creating the game; while reporting their game; playing their game; after editing their game), and for each time the participant resets the environment within each segment. We end up with 382 total such traces. Our primary use for them is in implementing a “reward machine,” an interpreter for our goal programs, which parses a goal program into a state machine, and iterates through a trace to emit the score of that trace under the goal. We use a limited version of this in our fitness features (see Fitness function methods for additional details) and in some of our model evaluations and ablations (see Supplementary information F for additional details).

Game dataset analyses methods

Common sense through predicate role-filler analysis: We analyze predicate role-filler occurrences, coarsening individual objects to higher-level categories (see the legend on the right of Figure ED-2). To split between the two panels of Figure ED-2, we categorize each game by whether it includes the following motifs: throwing (e.g., balls into a bin), stacking (e.g., blocks in a building), organizing (e.g., placing objects in specified places), or other. We split the figure into games involving only throwing motifs (left panel) and games involving any other motifs, potentially in addition to throwing (right panel). In games involving only throwing (left panel), participants most often refer to balls, primarily checking whether or not the agent holds a ball or a ball is in motion (as part of quantifying the act of throwing). Other predicates are often used to specify some additional conditions on throwing (such as specifying the bin being on the bed or the agent being next to the desk) and are used with a variety of object categories. Conversely, in games involving other elements (right panel), we see blocks and the generic “any_object” being used far more often, mostly in various placement and stacking constraints.

Compositionality and creativity through abstract structure occurrence: We analyze how often participant games make use of various grammatical structures to showcase both compositional reuse and long-tail creativity. Each structure involves a temporal modal (such as `once` or `hold`) and the

predicate expression nested under it, such as `(once (agent_holds ?b))`, where `?b` is a variable quantified earlier. We count structures, abstracting away specific variables and their types – so the expression above would be coarsened as `(once (agent_holds <obj>))`, and would be counted together with any other expression coarsened to this pattern. We encounter a total of 126 unique expressions in our dataset, the most common one with 62 occurrences being `(hold (and (not (agent_holds <obj>)) (in_motion <obj>)))`, which maps loosely to “find a sequence of states where an object is not held and is in motion” — that is, is currently moving with the agent touching it, for instance while being thrown or rolled. Of the 126 expressions, exactly half (63) occur only once.

Fitness function methods

Fitness function form: The fitness function used by our model is a learned, weighted linear combination of a set of features extracted programmatically from each game that is optimized to assign high scores to “human-like” games and low scores to everything else. It is a function $f : \mathcal{G} \rightarrow \mathbb{R}$ that maps individual games $g \in \mathcal{G}$ to real-valued scores: $f(g) = \theta \cdot \phi(g)$, where θ is a learned vector of weights and $\phi : \mathcal{G} \rightarrow [0, 1]^F$ is a feature extractor.

Feature extractor and feature set: The feature extractor ϕ represents each game as an 89-dimensional vector (i.e. $F = 89$). Each entry in the vector corresponds to a particular structural or semantic property of the game, from the size and depth of the syntax tree to the apparent feasibility of the game’s preferences. We normalize the values of each property to fall within the unit interval by using the observed range of values in our dataset. Many features used in the fitness function are directly computable from the DSL representation of a game (for instance, properties of its syntax tree or the misuse of particular grammatical structures). While these features represent the majority of the 89 features used, we also implement two important sets of features that require additional computation.

The first of these are n -gram features that capture the mean log score of the game under a simple n -gram language model trained over the set of human-generated syntax trees. We compute these scores separately for each game section (i.e. setup, preferences, terminal, and scoring) and also for the game overall, resulting in 5 features.

The second set consists of two features that make use of an interpreter that parses game programs into “reward machines” [34]: finite-state machines that process a trace of player inputs and emit a reward whenever the particular scoring conditions of the game are met. The interpreter programmatically implements each of the predicates in the DSL, which allows us to construct a dataset of which objects were used to satisfy which predicates across our dataset of 382 human *play traces*. The two features query this database in order to get an approximate common sense measure of a game’s “feasibility,” computing the proportion of a game’s predicate-argument combinations that have been satisfied by human players in our dataset (one feature does this for individual predicates, while the other does this for boolean logical expressions over predicates). While these feasibility measures give a sense of whether the objectives of a game can be completed in the physical reality of the simulation, the limited nature of our play trace dataset means they are far from perfect proxies.

The complete set of features used (and accompanying descriptions) is available in Supplementary information B, with the most important features (by their learned weights) highlighted in Supplementary information B.1.

Fitness function learning algorithm: To learn the weight vector θ , we take inspiration from the contrastive learning of energy-based models [11] with the objective of separating a set of *positive* examples (our set of human-generated games) from a set of *negative* examples. To learn an effective fitness function, these negatives must be qualitatively worse than our set of human games without being trivially distinguishable from them. We generate a set of plausible negatives by *corrupting* games from our positive set. To corrupt a game, we select a random node in its syntax tree, delete the node and its children, and randomly re-sample a sub-tree according to grammar of the DSL. This “tree-regrowth” approach [27] generally produces sub-trees that are syntactically valid but semantically “out-of-place,” with the severity of the corruption tending to correspond to height of the re-sampled node in the syntax tree. To account for the variance in the difficulty of distinguishing between a given positive and negative example, we generate a large set of negatives: 1024 for each of the 98 positives, for a total of 100,352 negatives.

We train the fitness function (i.e. optimize θ) using a softmax loss, not unlike the MEE loss used to train energy-based models [40] or the InfoNCE loss [67]. For a positive example g^+ and a set of negative examples $\{g_k^-\}, k \in \{1, 2, \dots, K\}$, we assign the loss:

$$\mathcal{L}(g^+, \{g_k^-\}_1^K; \theta) = -\log \frac{\exp(f_\theta(g^+))}{\exp(f_\theta(g^+)) + \sum_{k=1}^K \exp(f_\theta(g_k^-))} \quad (1)$$

This loss encourages the model to assign higher fitness scores to the real games than the negative examples. Simultaneously, this loss provides a diminishing incentive to push negative fitness scores down as the distance between the positives and negatives increases, intuitively assigning higher loss to negative examples with fitness closer to the positive example’s fitness. See Supplementary information C for full details of our training and cross-validation setups.

MAP-Elites methods

MAP-Elites overview: MAP-Elites is a population-based, evolutionary algorithm that works by defining a set of *behavioral characteristics*: discrete-valued functions over genotypes (in our case, game programs in the DSL) that form the axes of a multi-dimensional *archive* of cells. At each step, a game g is selected uniformly from among the individuals in the archive and mutated to form a new game g' . The mutated g' is evaluated both under the fitness function f and each of the n behavioral characteristics $b_i : \mathcal{G} \rightarrow \{0, \dots, k_i\}$ in order to determine which cell $c = [b_1(g), \dots, b_n(g)]$ it occupies. If the cell is unoccupied, then g' enters the archive. Otherwise, it enters the archive (and replaces the previous occupant) only if its fitness is greater than the current occupant of the cell. In this way, the algorithm maintains an “elite” for each possible combination of values under the behavioral characteristics.

Behavioral characteristics: Inspired by prior work on using MAP-Elites for procedural content generation [9], we define a set of integer-valued behavioral characteristics that each indicate how many preferences in each archive game match one of nine archetypal exemplar gameplay preferences. These include several types of ball-throwing preferences, as well as ones capturing block-stacking, object-sorting, and other miscellaneous activities. We also include two other features, one capturing whether or not the game includes a setup component, and one capturing the total number of preferences. For additional details and descriptions of the exemplar preferences, see Supplementary information D. The 11 total behavioral characteristics result in a total archive size of 2000 games.

Mutation operators: To mutate a game, we randomly select an operator from among the following: **regrowing** a random node and its children in its syntax tree, **inserting & deleting** the child of a node with multiple potential children, **crossing over** with the syntax tree of another randomly-selected game, **resampling the variables, initial conditions, or final conditions** used by a preference, and **resampling the optional game sections** (i.e. setup and terminal conditions). We seed the initial archive by naively sampling candidates from the PCFG—*not* with real, human-participant-created games or corruptions thereof that were used to train the fitness function. Further details of the algorithm are available in Supplementary information D.

Human evaluation methods

Evaluation dataset: we select games to be evaluated using the following procedure:

1. **real:** We include 30 participant-created games, each with a different set of behavioral characteristics — that is, each being considered ‘different’ according to how our model searches through the space of games (see MAP-Elites methods) for additional details).
2. **matched:** For each of the **real** games included above, we include the model-generated game from our final model from the corresponding MAP-Elites archive cell. Each of these games includes the same number of gameplay preferences as the corresponding **real** participant-created games, matching the same exemplar preferences.
3. **unmatched:** We then sample 30 additional games from our model’s archive. We sample in a fashion that aims to be balanced across the different preference counts and usage of the different exemplar preferences. That said, given that human games cover only 47 out of the 2000 archive cells,

that leaves 1953 potential unmatched games to sample; it is difficult to know how representative our set of 30 (which is about 1.5%) is.

GPT-4-based back-translation: Rather than ask participants to interpret our domain-specific language, we use the GPT-4 [52] language model to perform a multi-step back-translation from programs in our domain-specific language to structured natural language. For fairness and consistency, we use this procedure on the real games in addition to the model-generated matched and unmatched games. We first apply a rule-based system to apply templates, translating expressions in the DSL to natural language sentence fragments. We then use GPT-4 to first map the templated fragments to a more natural language, and then to combine the description of each game component (setup, gameplay preferences, terminal conditions, and scoring rules) to a short coherent description. See Supplementary information D for full details and prompts used.

Human evaluations structure: Figure ED-5 presents our human evaluation interface. Following instructions and an understanding quiz, participants evaluated nine total games: 3 real ones, the corresponding 3 matched ones, and 3 unmatched ones. Participants were presented one game at a time and provided two short textual responses, one explaining the game in their own words, and one providing a short overall impression of the game. Participants also answered seven Likert-type questions on 5-point scales, answering the following questions about the italicized attributes:

1. *Understandable*: “How confident are you that you understand the game described above?”, where 1: not at all confident, 3: moderately confident, and 5: very confident
2. *Fun to play*: “How fun would it be to play the game yourself?”, where 1: not at all fun, 3: moderately fun, and 5: extremely fun.
3. *Fun to watch*: “How fun would it be to watch someone else play this game?”, where 1: not at all fun, 3: moderately fun, and 5: extremely fun.
4. *Helpful*: “Imagine that you played this game for several minutes. How fun would it be for learning to interact with the virtual environment?”, where 1: not at all helpful, 3: moderately helpful, and 5: extremely helpful.
5. *Difficult*: “Imagine that you played this game for several minutes. Do you think it would be too easy, appropriately difficult, or too hard for you?”, where 1: far too easy, 3: appropriately difficult, and 5: far too hard.
6. *Creative*: “How creatively designed is this game?”, where 1: not at all creative, 3: moderately creative, and 5: extremely creative.
7. *Human-like*: “How human-like do you think this game is?”, where 1: not at all human-like, 3: moderately human-like, and 5: extremely human-like.

Evaluation statistical analyses: For each attribute and each game category (real, matched, and unmatched), we report the mean score assigned by all participants to games in that category for that attribute. We then also aggregate these attribute scores by category and report a nonparametric Mann-Whitney U -test [46] for differences in outcomes, as appropriate for ordinal data. See Supplementary Table SI-2 for the full table including test statistics and P-values. Significance results were highly similar when computing two-sample t -tests as an alternative statistical test.

Mixed effect models: We are interested in modeling the relationship between the scores predicted by our fitness function and the attributes human evaluators predicted. To that end, we set up mixed effect regression models [59, 31]. We fit separate models for each measure as the dependent variable, regressing a continuous latent score (e.g., s_{fp}^i for the fun-to-play measure, equation (2) below). We include fixed effects for the fitness score (x^i) and for membership in the real ($\mathbb{1}_{\text{real}}^i$) and matched ($\mathbb{1}_{\text{matched}}^i$) groups, treating the unmatched group as a baseline. We include random effects for the individual participants ($\epsilon_p^{pi} \sim \mathcal{N}(0, \sigma_p^2)$) and evaluated games ($\epsilon_g^{gi} \sim \mathcal{N}(0, \sigma_g^2)$). We also fit a sequence of cut-points (equation (3)) that transform the latent score to the observed ordinal rating y_{fp}^i (equation (4)). We suppress the subscript for each measure below:

$$s^i = \beta_{\text{fit}}x_i + \beta_{\text{real}}\mathbb{1}_{\text{real}}^i + \beta_{\text{matched}}\mathbb{1}_{\text{matched}}^i + \epsilon_p^{pi} + \epsilon_g^{gi} + \epsilon^i, \quad \epsilon^i \sim \mathcal{N}(0, \sigma^2) \quad (2)$$

$$-\infty \equiv c_0 < c_1 < c_2 < c_3 < c_4 < c_5 \equiv \infty \quad (3)$$

$$c_{k-1} < s^i \leq c_k \Rightarrow \text{observe } y^i = k \quad (4)$$

Models without either random effect performed worse than the full model, so we report results including both random effects. We fit cumulative link models for ordinal regression [2, 28] using the `ordinal` package [12] in **R** [58].

Sample similarity comparison methods

For each real game and its corresponding matched game from those included in the human evaluations, we examine which of our recorded participant interactions (see Dataset collection methods above) fulfills one or more gameplay elements. We treat the setup (if specified) and each gameplay preference as a gameplay element — our aim here is to quantify which participant interaction traces ‘play’ a part of the game. We do this using our “reward machine” — our implementation of an interpreter for goal programs in this domain-specific language. For each pair of games, we then check which particular interactions either (a) ‘play’ parts of both games, (b) only fulfill components in the real game, or (c) only fulfill components in the matched game. We color these proportions in purple, green, and blue respectively in Figure ED-6

TODOs

- Do we want to add another column with which features are active in each exemplar preference to the table below? 28

References

- [1] M. M. Andersen, J. Kiverstein, M. Miller, and A. Roepstorff. Play in predictive minds: A cognitive theory of play. *Psychological Review*, 130:462–479, 6 2023.
- [2] A. Argenti. *Categorical Data Analysis*. John Wiley & Sons, third edition, 2018.
- [3] J. T. Austin and J. B. Vancouver. Goal constructs in psychology: Structure, process, and content. *Psychological Bulletin*, 120:338–375, 11 1996.
- [4] D. E. Berlyne. Novelty and curiosity as determinants of exploratory behaviour1. *Br. J. Psychol. Gen. Sect.*, 41(1-2):68–80, Sept. 1950.
- [5] T. Brants, A. C. Papat, P. Xu, F. J. Och, and J. Dean. Large language models in machine translation. pages 858–867. Association for Computational Linguistics, 2007.
- [6] L. Brown. *Psychology of Motivation*. Nova Science Publishers, 2007.
- [7] C. Browne and F. Maire. Evolutionary game design. *IEEE Transactions on Computational Intelligence and AI in Games*, 2(1):1–16, 2010.
- [8] A. Campero, R. Raileanu, H. Küttler, J. B. Tenenbaum, T. Rocktäschel, and E. Grefenstette. Learning with amigo: Adversarially motivated intrinsic goals. 6 2021.
- [9] M. Charity, M. C. Green, A. Khalifa, and J. Togelius. Mech-elites: Illuminating the mechanic space of gvg-ai. In *Proceedings of the 15th International Conference on the Foundations of Digital Games*, pages 1–10, 2020.
- [10] K. Chatzilygeroudis, A. Cully, V. Vassiliades, and J. B. Mouret. Quality-diversity optimization: a novel branch of stochastic optimization. *Springer Optimization and Its Applications*, 170:109–135, 12 2020.
- [11] S. Chopra, R. Hadsell, and Y. LeCun. Learning a similarity metric discriminatively, with application to face verification. *Proceedings - 2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, CVPR 2005*, I:539–546, 2005.
- [12] R. H. B. Christensen. *ordinal—Regression Models for Ordinal Data*, 2023. R package version 2023.12-4.
- [13] J. Chu and L. E. Schulz. Play, curiosity, and cognition. *Annual Review of Developmental Psychology*, 2:317–343, 12 2020.
- [14] J. Chu, J. B. Tenenbaum, and L. E. Schulz. In praise of folly: flexible goals and human cognition. *Trends Cogn. Sci.*, Apr. 2024.
- [15] C. Colas, T. Karch, O. Sigaud, and P.-Y. Oudeyer. Autotelic agents with intrinsically motivated goal-conditioned reinforcement learning: a short survey. 12 2020.
- [16] C. Colas, L. Teodorescu, P.-Y. Oudeyer, X. Yuan, and M.-A. Côté. Augmenting autotelic agents with large language models. 5 2023.
- [17] A. Cully. Autonomous skill discovery with Quality-Diversity and unsupervised descriptors. May 2019.
- [18] Y. Du, O. Watkins, Z. Wang, C. Colas, T. Darrell, P. Abbeel, A. Gupta, and J. Andreas. Guiding pretraining in reinforcement learning with large language models. 7 2023.
- [19] C. S. Dweck. Article commentary: The study of goals in psychology. *Psychological Science*, 3(3):165–167, 1992.
- [20] J. S. Eccles and A. Wigfield. Motivational beliefs, values, and goals. *Annu. Rev. Psychol.*, 53:109–132, 2002.
- [21] A. J. Elliot and J. W. Fryer. The goal construct in psychology. In J. Y. Shah, editor, *Handbook of motivation science (pp*, volume 638, pages 235–250. The Guilford Press, xviii, New York, NY, US, 2008.
- [22] K. Fang, P. Yin, A. Nair, and S. Levine. Planning to practice: Efficient online fine-tuning by composing goals in latent space. 5 2022.
- [23] A. Fishbach and M. J. Ferguson. The goal construct in social psychology. In A. W. Kruglanski and E. T. Higgins, editors, *Social psychology: Handbook of basic principles*, volume 2, pages 490–515. The Guilford Press, xiii, New York, NY, US, 2007.

- [24] C. Florensa, D. Held, X. Geng, and P. Abbeel. Automatic goal generation for reinforcement learning agents. In *International conference on machine learning*, pages 1515–1528. PMLR, 2018.
- [25] J. A. Fodor. *The language of thought*. Harvard University Press, 1979.
- [26] M. Ghallab, A. Howe, C. Knoblock, D. McDermott, A. Ram, M. Veloso, D. Weld, and D. Wilkins. Pddl – the planning domain definition language. 1998.
- [27] N. D. Goodman, J. B. Tenenbaum, J. Feldman, and T. L. Griffiths. A rational analysis of rule-based concept learning. *Cogn. Sci.*, 32(1):108–154, Jan. 2008.
- [28] W. H. Greene and D. A. Hensher. *Modeling Ordered Choices: A Primer*. Cambridge University Press, 2010.
- [29] L. Grillotti and A. Cully. Unsupervised behaviour discovery with Quality-Diversity optimisation. June 2021.
- [30] V. Hom and J. Marks. Automatic design of balanced board games. In *Proceedings of the AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment*, volume 3, pages 25–30, 2007.
- [31] J. Hox, M. Moerbeek, and R. van de Schoot. *Multilevel Analysis (Techniques and Applications)*. New York, NY: Routledge, third edition, 2018.
- [32] W. Huang, F. Xia, T. Xiao, H. Chan, J. Liang, P. Florence, A. Zeng, J. Tompson, I. Mordatch, Y. Chebotar, et al. Inner monologue: Embodied reasoning through planning with language models. In *Conference on Robot Learning*, pages 1769–1782. PMLR, 2023.
- [33] M. E. Hyland. Motivational control theory: An integrative framework. *Journal of Personality and Social Psychology*, 55(4):642, 1988.
- [34] R. T. Icarte, T. Q. Klassen, R. Valenzano, and S. A. McIlraith. Using reward machines for high-level task specification and decomposition in reinforcement learning. 2018.
- [35] R. T. Icarte, T. Q. Klassen, R. Valenzano, and S. A. McIlraith. Reward machines: Exploiting reward function structure in reinforcement learning. *Journal of Artificial Intelligence Research* 73 (2022), 73:173–208, 10 2022.
- [36] A. Khalifa, M. C. Green, D. Perez-Liebana, and J. Togelius. General video game rule generation. In *2017 IEEE Conference on Computational Intelligence and Games (CIG)*, pages 170–177. IEEE, 2017.
- [37] E. Kolve, R. Mottaghi, W. Han, E. VanderBilt, L. Weihs, A. Herrasti, M. Deitke, K. Ehsani, D. Gordon, Y. Zhu, et al. Ai2-thor: An interactive 3d environment for visual ai. *arXiv preprint arXiv:1712.05474*, 2017.
- [38] B. M. Lake, T. D. Ullman, J. B. Tenenbaum, and S. J. Gershman. Building machines that learn and think like people. *Behav. Brain Sci.*, 40:e253, Jan. 2017.
- [39] P. H. Le-Khac, G. Healy, and A. F. Smeaton. Contrastive representation learning: A framework and review. Oct. 2020.
- [40] Y. LeCun, S. Chopra, R. Hadsell, M. Ranzato, and F. J. Huang. *A Tutorial on Energy-Based Learning*. MIT Press, 2006.
- [41] R. V. Lenth. *emmeans: Estimated Marginal Means, aka Least-Squares Means*, 2024. R package version 1.10.0.
- [42] B. G. Leon, M. Shanahan, and F. Belardinelli. In a nutshell, the human asked for this: Latent goals for following temporal specifications openreview. *ICLR 2022*.
- [43] A. S. Lillard. *The Development of Play*, volume 3, pages 425–468. Wiley-Blackwell, 2015.
- [44] M. L. Littman, U. Topcu, J. Fu, C. Isbell, M. Wen, and J. MacGlashan. Environment-independent task specifications via gltl. *arXiv*, 4 2017.
- [45] M. Liu, M. Zhu, and W. Zhang. Goal-Conditioned reinforcement learning: Problems and solutions. Jan. 2022.
- [46] H. B. Mann and D. R. Whitney. On a test of whether one of two random variables is stochastically larger than the other. *Ann. Math. Stat.*, 18(1):50–60, 1947.

- [47] D. Marr. *Vision: A Computational Investigation into the Human Representation and Processing of Visual Information*. Henry Holt and Co., Inc., USA, 1982.
- [48] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.
- [49] G. Molinaro and A. G. E. Collins. A goal-centric outlook on learning. *Trends Cogn. Sci.*, 27(12):1150–1164, Dec. 2023.
- [50] G. B. Moskowitz and H. Grant, editors. *The psychology of goals*, volume 548. Guilford Press, New York, NY, US, 2009.
- [51] J.-B. Mouret and J. Clune. Illuminating search spaces by mapping elites. *arXiv preprint arXiv:1504.04909*, 2015.
- [52] OpenAI. GPT-4 technical report. Mar. 2023.
- [53] P.-Y. Oudeyer, F. Kaplan, and V. V. Hafner. Intrinsic motivation systems for autonomous mental development. *IEEE Trans. Evol. Comput.*, 11(2):265–286, Apr. 2007.
- [54] S. Palan and C. Schitter. Prolific.ac—A subject pool for online experiments. *Journal of Behavioral and Experimental Finance*, 17:22–27, Mar. 2018.
- [55] B. Pell. Metagame in symmetric chess-like games. 1992.
- [56] L. Pervin. *Goal Concepts in Personality and Social Psychology*. Psychology Library Editions: Social Psychology. Taylor & Francis, 2015.
- [57] J. K. Pugh, L. B. Soros, and K. O. Stanley. Quality diversity: A new frontier for evolutionary computation. *Frontiers in Robotics and AI*, 3, 2016.
- [58] R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2023.
- [59] S. W. Radenbush and A. S. Bryk. *Hierarchical Linear Models. Applications and Data Analysis Methods*. Thousand Oaks, CA: Sage Publications, second edition, 2002.
- [60] A. Ruggeri, O. Stanciu, M. Pelz, A. Gopnik, and E. Schulz. Preschoolers search longer when there is more information to be gained. *Dev. Sci.*, 27(1):e13411, Jan. 2024.
- [61] J. S. Rule, J. B. Tenenbaum, and S. T. Piantadosi. The Child as Hacker. *Trends in Cognitive Sciences*, 24(11):900–915, nov 2020.
- [62] A. M. Smith, M. J. Nelson, and M. Mateas. Ludocore: A logical game engine for modeling videogames. In *Proceedings of the 2010 IEEE Conference on Computational Intelligence and Games*, pages 91–98. IEEE, 2010.
- [63] R. S. Sutton and A. G. Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- [64] O. E. L. Team, A. Stooke, A. Mahajan, C. Barros, C. Deck, J. Bauer, J. Sygnowski, M. Trebacz, M. Jaderberg, M. Mathieu, et al. Open-ended learning leads to generally capable agents. *arXiv preprint arXiv:2107.12808*, 2021.
- [65] A. Ten, P. Kaushik, P.-Y. Oudeyer, and J. Gottlieb. Humans monitor learning progress in curiosity-driven exploration. *Nat. Commun.*, 12(1):5972, Oct. 2021.
- [66] J. Togelius and J. Schmidhuber. An experiment in automatic game design. In *2008 IEEE Symposium On Computational Intelligence and Games*, pages 111–118. IEEE, 2008.
- [67] van den Oord Aaron, Y. Li, and O. Vinyals. Representation learning with contrastive predictive coding. 7 2018.
- [68] G. Wang, Y. Xie, Y. Jiang, A. Mandlekar, C. Xiao, Y. Zhu, L. Fan, and A. Anandkumar. Voyager: An open-ended embodied agent with large language models. *arXiv preprint arXiv:2305.16291*, 2023.
- [69] T. B. Ward. Structured imagination: the role of category structure in exemplar generation. *Cognitive Psychology*, 27:1–40, 8 1994.

- [70] L. Wong, G. Grand, A. K. Lew, N. D. Goodman, V. K. Mansinghka, J. Andreas, and J. B. Tenenbaum. From word models to world models: Translating from natural language to the probabilistic language of thought. June 2023.
- [71] A. Zook and M. Riedl. Automatic game design via mechanic generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 28, 2014.

Extended Data

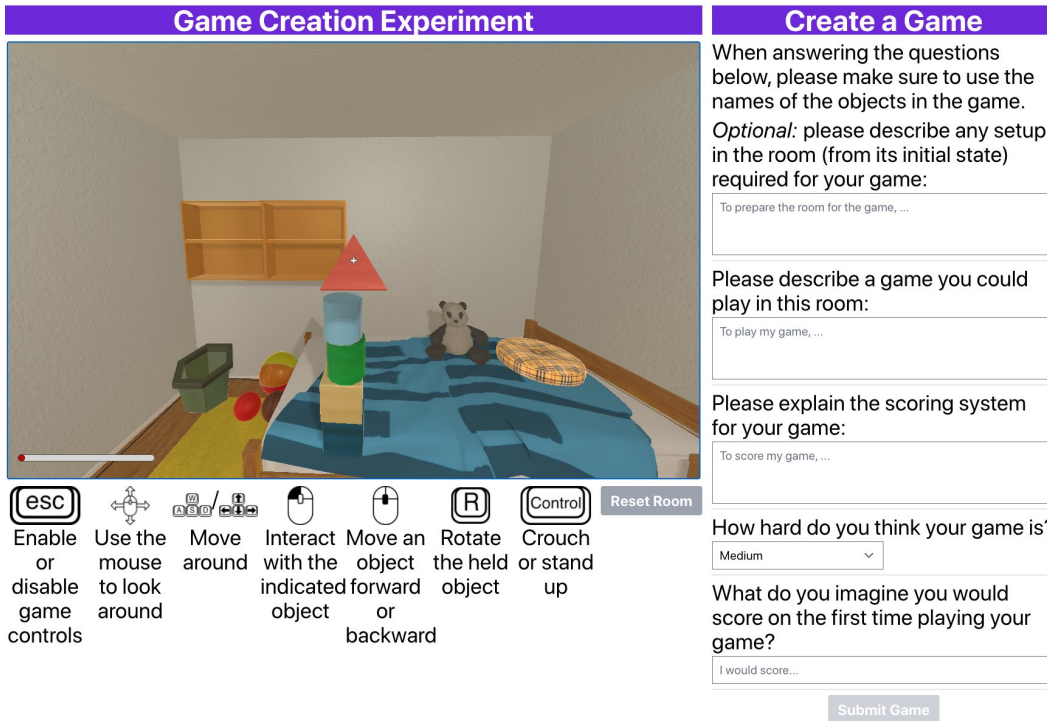


Figure ED-1: Online experiment interface. The main part of the screen presents the AI2-THOR-based experiment room. Below it, we depict the controls. To the right, we show the text prompts for creating a new game (fonts enlarged for visualization). Our experiment is accessible online here.

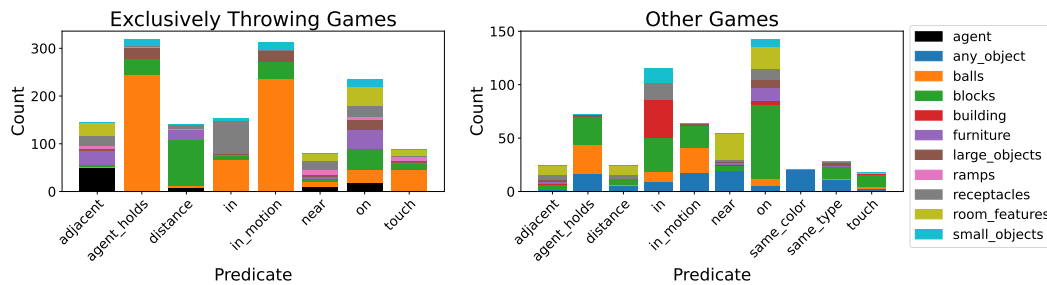


Figure ED-2: Participants showcase intuitive common sense. **Left:** In games involving exclusively throwing, participants use balls (orange) far more often than any other object type. **Right:** In other games, participants refer to blocks or “any object” more often, most often checking where objects are placed (using the **in** and **on** predicates). We most often observe balls being thrown and blocks being stacked, and while a few participants specified block-throwing games, no participant created a game involving ball-stacking. Other types of common sense were also recruited: participants rarely specified throwing games involving large or cumbersome objects (such as the chair or laptop), and mostly specified stacking objects in objects in buildings (as opposed to trying to move buildings or touch them). See Game dataset analyses methods for additional details.

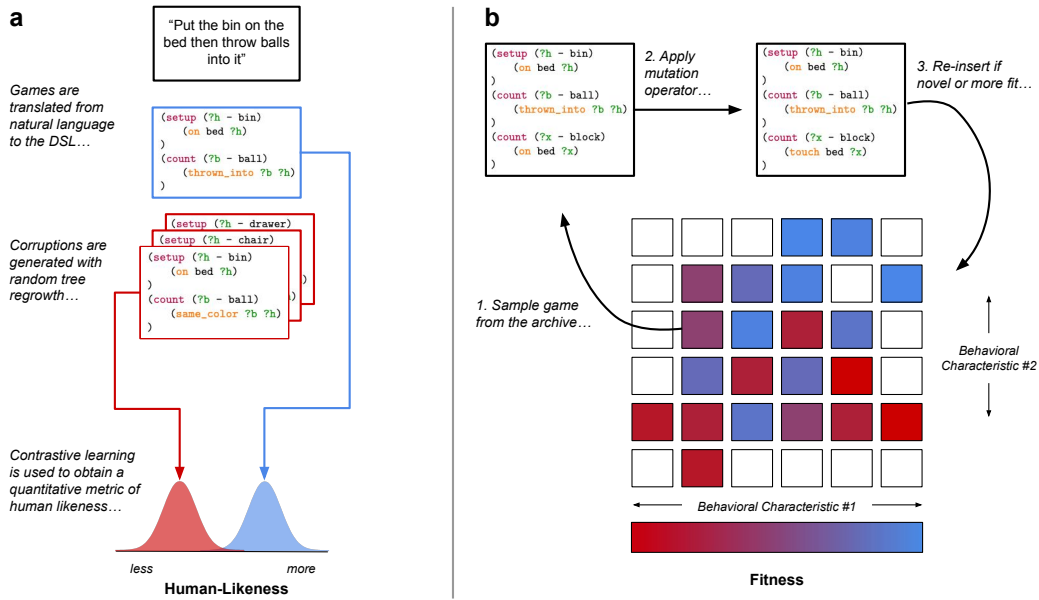


Figure ED-3: Model architecture. **Left:** We contrastively learn a quantitative measure of human likeness by maximizing the distance between human-generated exemplar games and a set of corruptions obtained through random tree regrowth. **Right:** This measure is then used as the basis for quality-diversity optimization through MAP-Elites. The algorithm maintains an archive of games that differ across phenotypic “behavioral characteristics.” At each step, a game is randomly sampled from the archive (1), randomly mutated (2), and re-evaluated for both fitness and its position in the archive. It is added to the archive only if it would occupy a previously empty position in the archive or if it is more fit than the current occupant (3).

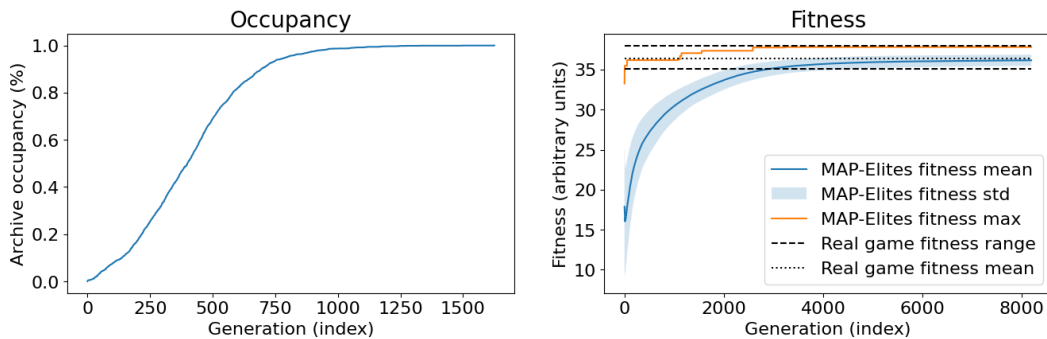
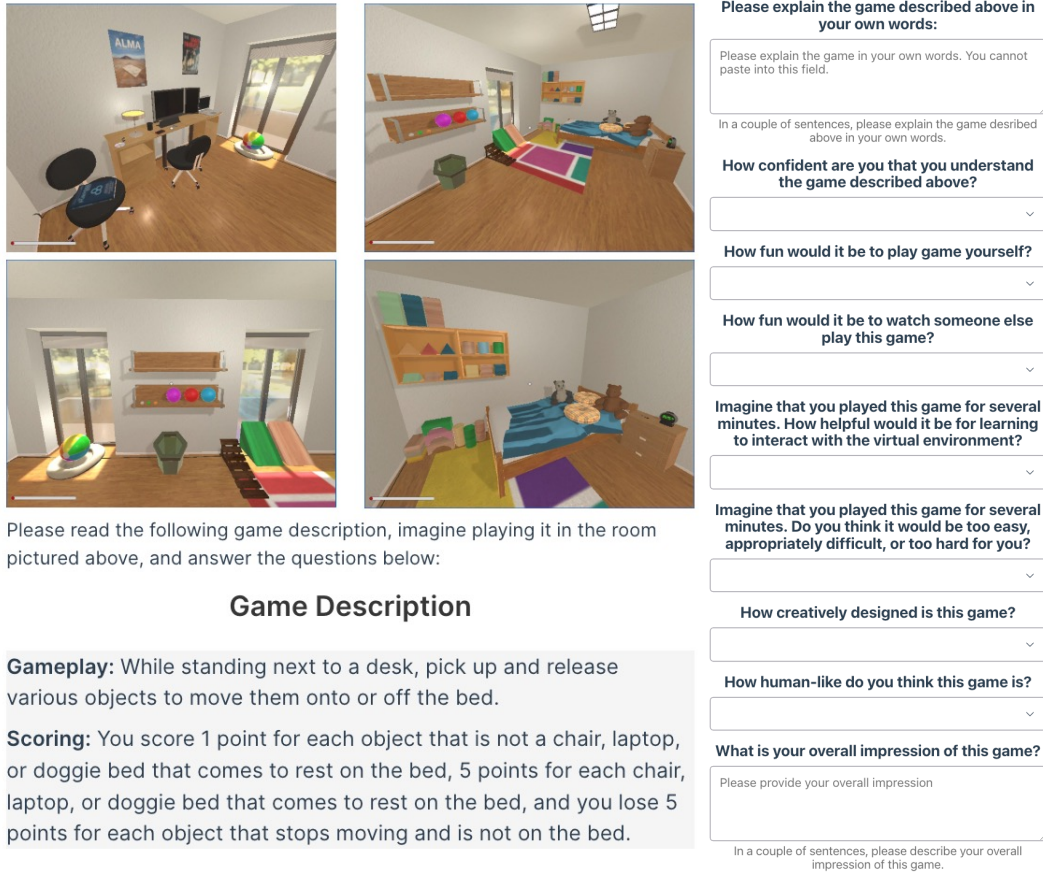


Figure ED-4: Our model fills the archive quickly and finds examples with human-like fitness scores. **Left:** Our model rapidly finds exemplars for all archive cells (i.e. niches induced by our behavioral characteristics), reaching 50% occupancy after 400 generations (out of a total of 8192) and 95% occupancy after 794 generations—the archive is almost full 1/10th of the way through the search process. **Right:** Our model reaches human-like fitness scores. After only three generations, the fittest sample in the archive has a higher fitness score than at least one participant-created game. By the end of the search, the mean fitness in the archive is close to the mean fitness of human games.



Please read the following game description, imagine playing it in the room pictured above, and answer the questions below:

Game Description

Gameplay: While standing next to a desk, pick up and release various objects to move them onto or off the bed.

Scoring: You score 1 point for each object that is not a chair, laptop, or doggie bed that comes to rest on the bed, 5 points for each chair, laptop, or doggie bed that comes to rest on the bed, and you lose 5 points for each object that stops moving and is not on the bed.

Please explain the game described above in your own words:

Please explain the game in your own words. You cannot paste into this field.

In a couple of sentences, please explain the game described above in your own words.

How confident are you that you understand the game described above?

How fun would it be to play game yourself?

How fun would it be to watch someone else play this game?

Imagine that you played this game for several minutes. How helpful would it be for learning to interact with the virtual environment?

Imagine that you played this game for several minutes. Do you think it would be too easy, appropriately difficult, or too hard for you?

How creatively designed is this game?

How human-like do you think this game is?

What is your overall impression of this game?

Please provide your overall impression

In a couple of sentences, please describe your overall impression of this game.

Figure ED-5: Human evaluations interface. For each game, participants viewed the same four images of the environment, followed by the GPT-4 back-translated description of the game (see Human evaluation methods for details). They then answered the two free-response and seven multiple-choice questions on the right. In the web-page based version, the questions appeared below the game description; they are presented side-by-side to save space.

Table ED-1: Mixed model result summary

| Measure | Fitness | | Variable | | ℓ [Real] | |
|-------------------------|---------------|--------------|-------------------|--------------|----------------|--------------|
| | β_{fit} | Significance | $\beta_{matched}$ | Significance | β_{real} | Significance |
| <i>Understandable</i> ↑ | 0.846 | *** | 0.525 | - | 1.151 | *** |
| <i>Fun to play</i> ↑ | 0.396 | ** | 0.629 | * | 1.059 | *** |
| <i>Fun to watch</i> ↑ | 0.191 | - | 0.641 | * | 0.912 | *** |
| <i>Helpful</i> † ↑ | -0.189 | * | 0.349 | * | 0.232 | - |
| <i>Difficult</i> ↓ | -0.588 | *** | 0.363 | - | -0.250 | - |
| <i>Creative</i> ↑ | -0.486 | ** | 0.551 | - | 0.438 | - |
| <i>Human-like</i> ↑ | 0.570 | *** | 0.837 | ** | 1.446 | *** |

Fitness scores significantly predict several attributes, including understandability and human-likeness. Fitness scores show (statistically) significant positive effects on the understandability, fun to play, and human-likeness attributes, and significant negative effects on the helpfulness, difficulty and creativity questions. Accounting for the role of fitness, the *matched* group membership shows significant effects only the fun to play and watch, helpfulness, and human likeness questions. The *real* group shows significant effects on understandability, fun to play and watch, and human likeness. See Supplementary Table SI-3 for test statistics and P-values. *: $P < 0.05$, **: $P < 0.01$, ***: $P < 0.001$ †: The full measure description is “Helpful for interacting with the simulated environment.”

In most measures, higher scores are better, indicated by the ↑, other than *Difficult* ↓, in which 3 means “appropriately difficult”, and scores below and above indicate too easy and too hard respectively.

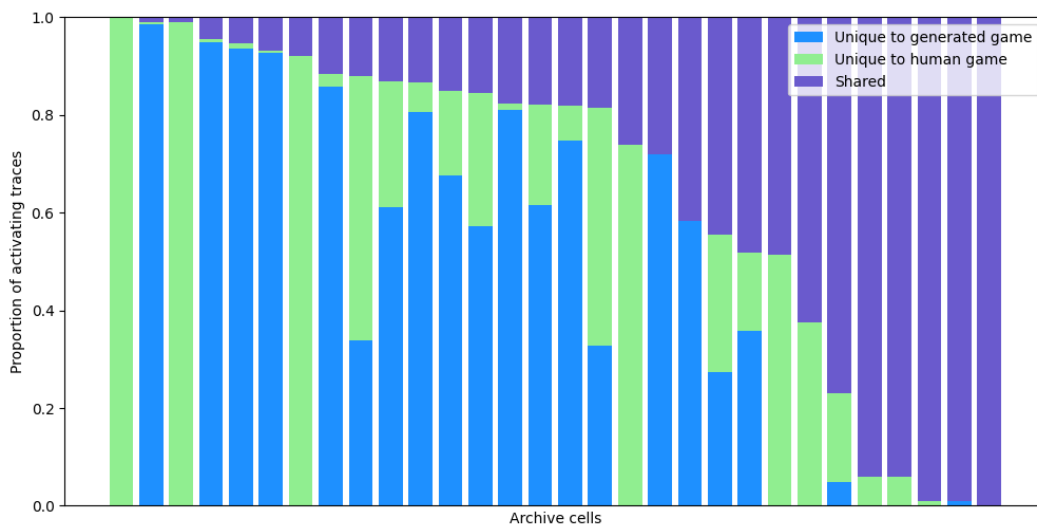


Figure ED-6: Proportion of human interactions activating only matched and real games in the same cell. Each bar corresponds to a pair of corresponding matched and real games. In each bar, we plot the proportion of relevant interactions (state-action traces) that are unique to the matched game (blue), unique to the real game (green), or shared across both (purple). A few games (with the bar mostly or entirely in purple) show high similarity between the corresponding games — under 25% (7/30) share more than half of their relevant interactions. Most games, however, show substantial differences between the sets of relevant interactions, with some showing a higher fraction unique to human games and others to matched model games. The average Jaccard similarity between the sets of relevant interactions for the matched and real game is only 0.347 and the median similarity is 0.180 (identical games would score 1.0, entirely dissimilar games 0).

B Full Feature Set

To simplify training fitness models, we ensure that all feature values are on the unit interval, using the following feature types:

- A binary value (marked with [b])
- A proportion between zero and one ([p])
- A real value discretized to two or more levels and treated as an indicator variable ([d], with the levels listed at the end of the description)
- A float value normalized to the unit interval over the full dataset of positive and negative games ([f])

For our n-gram features, we extract n-gram tokens from an in-order traversal of the syntax tree. We use 5-gram models with stupid backoff [5] with a discount factor of 0.4, and report the mean log score as the feature value, both jointly over the entire game program and separately over the different sections (setup, preferences, terminal conditions, and scoring).

For our predicate play trace features, we use a simplified version of the predicate satisfaction computation aspect of our reward machine (DSL program interpreter). We record, for every human play trace we have, and each predicate listed below, for every object assignment that satisfies it in that trace, all indices of states at which the predicate is satisfied. Recording specific states allows to us compute conjunctions, disjunctions, and negations in addition to individual predicate satisfactions. We limit ourselves to a subset of our predicates, which covers over 95% of predicate references in our dataset: `above`, `adjacent`, `agent_crouches`, `agent_holds`, `broken`, `game_start`, `game_over`, `in`, `in_motion`, `object_orientation`, `on`, `open`, `toggled_on`, and `touch`. Any predicate that is not implemented is assumed to be feasible to have been satisfied.

Our full feature set is:

ngram: Features using our n-gram model.

1. `ast_ngram_full_n_5_score` [f]: What is the mean 5-gram model score under an n-gram model trained on the real games?
2. `ast_ngram_setup_n_5_score` [f]: What is the mean 5-gram model score of the setup section under an n-gram model trained on the real game setup sections?
3. `ast_ngram_constraints_n_5_score` [f]: What is the mean 5-gram model score of the gameplay preferences section under an n-gram model trained on the real game preferences sections?
4. `ast_ngram_terminal_n_5_score` [f]: What is the mean 5-gram model score of the terminal conditions section under an n-gram model trained on the real game terminal sections?
5. `ast_ngram_scoring_n_5_score` [f]: What is the mean 5-gram model score of the scoring section under an n-gram model trained on the real game scoring sections?

play_trace_database: Features using our play trace database.

1. `predicate_found_in_data_prop` [p]: What proportion of predicates are satisfied at least once in our human play trace data?
2. `predicate_found_in_data_small_logicals_prop` [p]: What proportion of logical expressions over predicates (with four or fewer children, limited for computational reasons) are satisfied at least once in our human play trace data?

defined_and_used: Features reflecting whether particular game components are defined, and features capturing whether defined components (such as variables, gameplay preferences, or objects in the setup) are then also used elsewhere.

1. `variables_used_all` [b]: Are all variables defined used at least once?
2. `variables_used_prop` [p]: What proportion of variables defined are used at least once?
3. `preferences_used_all` [b]: Are all preferences defined referenced at least once in terminal or scoring expressions?
4. `preferences_used_prop` [p]: What proportion of preferences defined are referenced at least once in terminal or scoring expressions?

5. `setup_quantified_objects_used` [p]: What proportion of object or types quantified as variables in the setup are also referenced in the gameplay preferences?
6. `any_setup_objects_used` [b]: Are any objects referenced in the setup also referenced in the gameplay preferences?
7. `section_doesnt_exist_setup` [b]: Does a game not have an (optional) setup section? (to allow counteracting feature values for the setup for games that do not have a setup component)
8. `section_doesnt_exist_terminal` [b]: Does a game not have an (optional) terminal conditions section? (to allow counteracting feature values for the terminal conditions for games that do not have a terminal conditions component)

grammar_misuse: Features capturing various modes of grammar misuse—expressions that are grammatical under the DSL but ill-formed, poorly structured, or whose values cannot vary over gameplay.

1. `adjacent_once_found` [b]: Are there any cases where the `once` modal, which captures a single state, is used twice in a row?
2. `adjacent_same_modal_found` [b]: Are there any cases where the same modal is used twice in a row?
3. `once_in_middle_of_pref_found` [b]: Are there any cases where the `once` modal, which captures a single state, is in the middle of a sequence of modals?
4. `pref_without_hold_found` [b]: Are there any cases where a sequence of modals is specified with no temporally extended modal (`hold` or `hold-while`)?
5. `identical_consecutive_seq_func_predicates_found` [b]: Are there any cases where the same exact predicates (and their arguments) are applied in consecutive modals (making them redundant)?
6. `predicate_without_variables_or_agent` [b]: Are there any predicates that do not reference any variables or the agent?
7. `nested_logicals_found` [b]: Are there any cases where a logical operator is nested inside the same logical operator (e.g., a negation of a negation, or a conjunction of a conjunction)?
8. `identical_logical_children_found` [b]: Are there any cases where a logical operator has two or more identical children?
9. `redundant_expression_found` [b]: Are there any cases where a logical expression over predicates is redundant (can be trivially simplified)?
10. `unnecessary_expression_found` [b]: Are there any cases where a logical expression over predicates is unnecessary (contradicts itself, or is trivially true)?
11. `repeated_variables_found` [b]: Are there any cases where the same variable is used twice in the same predicate?
12. `repeated_variable_type_in_either` [b]: Are there any cases where the same variable types is used twice in an `either` quantification?

scoring_grammar_misuse: Features capturing similar failure modes to the above category, but localized to the scoring and terminal sections of the DSL.

1. `identical_scoring_children_found` [b]: Are there any cases where a scoring arithmetic or logical expression has two or more identical children?
2. `redundant_scoring_terminal_expression_found` [b]: Are there any cases where a scoring or terminal expression is redundant (can be trivially simplified)?
3. `unnecessary_scoring_terminal_expression_found` [b]: Are there any cases where a scoring or terminal expression is unnecessary (contradicts itself, or is trivially true)?
4. `total_score_non_positive` [b]: Do the scoring rules of the game result in a non-positive score regardless of gameplay?
5. `scoring_preferences_used_identically` [b]: Do the scoring rules of the game treat all gameplay preferences identically?
6. `two_number_operation_found` [b]: Are there any cases where an arithmetic operation is applied to two numbers? (e.g. `(+ 5 5)` instead of simplifying it)

game_element_disjointness: Features capturing whether particular game elements are disjoint—for example, gameplay preferences using disjoint sets of objects, or temporal modals using disjoint sets of variables.

1. `disjoint_preferences_found` [b]: Are there any preferences that quantify over disjoint sets of objects?

2. `disjoint_preferences_scoring_terminal_types` [p]: Do the preferences referenced in the scoring and terminal section quantify over disjoint sets of object types?
3. `disjoint_preferences_scoring_terminal_predicates` [p]: Do the preferences referenced in the scoring and terminal section use disjoint sets of predicates?
4. `disjoint_seq_funcs_found` [b]: Are there any cases where modals in a preference refer to disjoint sets of variables or objects?
5. `disjoint_at_end_found` [b]: Are there any cases where predicate expressions under an `at_end` refer to disjoint sets of variables or objects?
6. `disjoint_modal_predicates_found` [b]: Are there any cases where modals in a preference refer to disjoint sets of predicates?
7. `disjoint_modal_predicates_prop` [p]: What proportion of modals in a preference refer to disjoint sets of predicates?

counting: Features tracking node count or maximal depth in the four different DSL program sections.

1. `node_count_section` [d]: How many nodes are in the *section*, discretized to five bins with different thresholds for each section.
2. `max_depth_section` [d]: What is the maximal depth of the syntax tree in the *section*, discretized to five bins with different thresholds for each section.

pref_forall: Features capturing whether or not and how well the games use the `forall` over preferences syntax.

1. `pref_forall_used_correct` [b]: For the `forall` over preferences syntax, if it is used, is it used correctly in this game?
2. `pref_forall_used_incorrect` [b]: For the `forall` over preferences syntax, if it is used, is it used incorrectly in this game? (to allow learning differential values between correct and incorrect usage)
3. `pref_forall_external_forall_used_correct` [b]: If the `count-once-per-external-objects` count operator is used, is it used correctly in this game?
4. `pref_forall_external_forall_used_incorrect` [b]: If the `count-once-per-external-objects` count operator is used, is it used incorrectly in this game?
5. `pref_forall_pref_forall_correct_arity_correct` [b]: If optional object names and types are provided to a count operation, are they provided with an arity consistent with the `forall` variable quantifications?
6. `pref_forall_pref_forall_correct_arity_incorrect` [b]: If optional object names and types are provided to a count operation, are they provided with an arity inconsistent with the `forall` variable quantifications?
7. `pref_forall_pref_forall_correct_types_correct` [b]: If optional object names and types are provided to a count operation, are they provided with types consistent with the `forall` variable quantifications?
8. `pref_forall_pref_forall_correct_types_incorrect` [b]: If optional object names and types are provided to a count operation, are they provided with types inconsistent with the `forall` variable quantifications?

B.1 Features Most Predictive of Real or Regrown Games

The following features (in order) had the largest weight, indicating they were most predictive of positive (real, human-generated) examples in our dataset. The last three features all capture the same concept, whether or not a setup section exists. We surmise the diffused weights over them are a result of using weight decay (an L2 penalty) on the model weights:

1. `ast_ngram_full_n_5_score`
2. `ast_ngram_constraints_n_5_score`
3. `predicate_found_in_data_prop`
4. `ast_ngram_setup_n_5_score`
5. `variables_used_all`
6. `preferences_used_all`
7. `ast_ngram_scoring_n_5_score`

8. `max_depth_setup_0` (which indicates a setup section does not exist or is very minimal)
9. `node_count_setup_0` (which indicates a setup section does not exist or is very minimal)
10. `section_doesnt_exist_setup`

The following features (in order) had the smallest weights, indicating they were most predictive of negative (regrown) examples in our dataset:

1. `pref_forall_used_incorrect`
2. `pref_forall_pref_forall_correct_types_incorrect`
3. `disjoint_seq_funcs_found`
4. `repeated_variables_found`
5. `redundant_expression_found`
6. `pref_forall_pref_forall_correct_arity_incorrect`
7. `predicate_without_variables_or_agent`
8. `two_number_operation_found`
9. `nested_logicals_found`
10. `redundant_scoring_terminal_expression_found`

C Objective function algorithm descriptions

Algorithm 1 below outlines how we train our fitness model. The number N of positive examples is fixed (98 in our full dataset), and fewer during cross-validation. We generate $M = 1024$ negatives for each of the positive examples, and the number of features F is fixed as well. We perform cross-validation to select hyperparameter values $B \in \{1, 2, 4\}$, and $K \in \{256, 512, 1025\}$, selecting the set that minimizes the cross-validated loss. We optimize the model with SGD, with a learning rate $\eta \in \{1e-3, 4e-3\}$ also selected via cross-validation. We use weight decay with $\lambda = 0.003$ to regularize the model. We train the model for up to 25000 epochs, or until the model plateaus for $P = 500$ epochs. After cross-validation, we train our final objective function on the entire dataset. The final model we report uses $B = 1$ positive games per batch, $K = 1024$ negative samples from our entire dataset for that positive, a learning rate $\eta = 4e-3$, and $F = 50$ features.

D MAP-Elites Algorithm Details

We use a set of 9 exemplar preferences as the basis for our MAP-Elites behavioral characteristics, detailed in Table SI-1. To score each game with respect to each exemplar preference, we count how many of the game's preferences are a close match to the exemplar. We explored matching preferences by edit distance (in string or syntax tree space) but discovered the edit distance is rather easily game-able by the model, producing semantically similar preferences with high edit distance from each other. Instead, we represent each exemplar preference as a binary feature vector, with features for which groups of predicates the preference uses (4 features: `agent_holds` or `in_motion`, `in`, `on`, and `adjacent` or `near` or `touch`), and for which object categories the preference uses (5 features: `balls`, `receptacles`, `blocks` or `buildings`, `furniture` or `room_features`, and `small_items` or `large_items` or the generic `game_object`). Preferences in each archive game are also represented using this feature space. A preference in an archive game is considered a match for an exemplar if it has an L1 distance of 0 or 1 in this space, and if a preference matches more than one exemplar, a match is randomly chosen. Exemplar preferences were iteratively chosen, starting from a seed preference (the first in Table SI-1), and then greedily adding additional exemplars from the preferences defined in participant-created games. At each step, the preference added was chosen to maximize the number of participant-created preferences that would be considered a match (distance of 0 or 1) from the exemplar set. In addition, we include two other behavioral characteristics, one capturing whether or not the game includes a setup component, and one capturing the total number of

Do we want to add another column with which features are active in each exemplar preference to the table below?

Table SI-1: Exemplar preferences used as MAP-Elites behavioral characteristics.

| Exemplar Preference | Description (GPT-4 back-translated) | Exemplar Features |
|--|---|--|
| <pre>(preference throwAttempt (exists (?b - dodgeball) (then (once (agent_holds ?b)) (hold (and (not (agent_holds ?b)) (in_motion ? b))) (once (not (in_motion ?b)))))) </pre> | <p>This preference is satisfied when:</p> <ul style="list-style-type: none"> -first, the agent holds a dodgeball -next, the agent throws the dodgeball -finally, the dodgeball stops moving | <p>Uses predicate <code>agent_holds</code> or <code>in_motion</code></p> <p>Uses object category <code>balls</code></p> |
| <pre>(preference throwInBin (exists (?b - ball ?h - hexagonal_bin) (then (once (and (on rug agent) (agent_holds ?b))) (hold (and (not (agent_holds ?b)) (in_motion ? b))) (once (and (not (in_motion ?b)) (in ?h ?b)))))) </pre> | <p>This preference is satisfied when:</p> <ul style="list-style-type: none"> -first, the agent is standing on the rug and holding a ball -next, the agent throws the ball -finally, the ball stops moving and is inside a hexagonal bin | <p>Uses predicate <code>agent_holds</code> or <code>in_motion</code></p> <p>Uses predicate <code>in</code></p> <p>Uses predicate <code>on</code></p> <p>Uses object category <code>balls</code></p> <p>Uses object category <code>receptacles</code></p> <p>Uses object category <code>furniture</code> or <code>room_features</code></p> |
| <pre>(preference ballThrownToBed (exists (?d - dodgeball) (then (once (and (agent_holds ?d) (adjacent desk agent))) (hold (and (not (agent_holds ?d)) (in_motion ? d))) (once (and (not (in_motion ?d)) (on bed ?d)))))) </pre> | <p>This preference is satisfied when:</p> <ul style="list-style-type: none"> -first, the agent holds a dodgeball while standing next to a desk -next, the agent throws the dodgeball -finally, the dodgeball stops moving and is on the bed | <p>Uses predicate <code>agent_holds</code> or <code>in_motion</code></p> <p>Uses predicate <code>on</code></p> <p>Uses predicate <code>adjacent</code> or <code>near</code> or <code>touch</code></p> <p>Uses object category <code>balls</code></p> <p>Uses object category <code>furniture</code> or <code>room_features</code></p> |
| <pre>(preference itemInClosedDrawerAtEnd (exists (?g - game_object) (at_end (and (in top_drawer ?g) (not (open top_drawer))))))) </pre> | <p>This preference is satisfied when:</p> <ul style="list-style-type: none"> -at the end of the game, a game object is inside the top drawer and the top drawer is closed | <p>Uses predicate <code>in</code></p> <p>Uses object category <code>receptacles</code></p> <p>Uses object category <code>small_objects</code> or <code>large_objects</code> or <code>any_object</code></p> <p>Uses <code>at_end</code></p> |
| <pre>(preference watchOnShelf (exists (?w - watch ?s - shelf) (at_end (on ?s ?w)))) </pre> | <p>This preference is satisfied when:</p> <ul style="list-style-type: none"> -at the end of the game, a watch is on a shelf | <p>Uses predicate <code>on</code></p> <p>Uses object category <code>furniture</code> or <code>room_features</code></p> <p>Uses object category <code>small_objects</code> or <code>large_objects</code> or <code>any_object</code></p> <p>Uses <code>at_end</code></p> |
| <pre>(preference gameBlockFound (exists (?1 - block) (then (once (game_start)) (hold (not (exists (?b - building) (and (in ? b ?1) (is_setup_object ?b)))))) (once (agent_holds ?1))))) </pre> | <p>This preference is satisfied when:</p> <ul style="list-style-type: none"> -first, the game begins -next, throughout the game, the block is not part of a building that is used in the setup -finally, the agent picks up the block | <p>Uses predicate <code>agent_holds</code> or <code>in_motion</code></p> <p>Uses predicate <code>in</code></p> <p>Uses object category <code>blocks</code> or <code>building</code></p> |
| <pre>(preference matchingBuildingBuilt (exists (?b1 ?b2 - building) (at_end (and (is_setup_object ?b1) (not (is_setup_object ?b2)) (forall (?1 ?12 - block) (or (not (in ?b1 ?11)) (not (in ?b1 ?12)) (not (on ?11 ?12)) (exists (?13 ?14 - block) (and (in ?b2 ?13) (in ?b2 ?14) (on ?13 ?14) (same_type ?11 ?13) (same_type ?12 ?14))))))))))))) </pre> | <p>This preference is satisfied when:</p> <ul style="list-style-type: none"> -at the end of the game, one building is part of the setup while the other is not -and for any two blocks, neither is inside the building that is part of the setup -if one block is not on top of the other, then there must be two other blocks of the same types inside the building that is not part of the setup, with one of these blocks on top of the other | <p>Uses predicate <code>in</code></p> <p>Uses predicate <code>on</code></p> <p>Uses object category <code>blocks</code> or <code>building</code></p> <p>Uses <code>at_end</code></p> |
| <pre>(preference ballDroppedInBin (exists (?b - ball ?h - hexagonal_bin) (then (once (and (adjacent ?h agent) (agent_holds ? b))) (hold (and (in_motion ?b) (not (agent_holds ? b)))) (once (and (not (in_motion ?b)) (in ?h ?b)))))) </pre> | <p>This preference is satisfied when:</p> <ul style="list-style-type: none"> -first, the agent is next to a hexagonal bin and is holding a ball -next, the agent throws the ball -finally, the ball stops moving and is inside the hexagonal bin | <p>Uses predicate <code>agent_holds</code> or <code>in_motion</code></p> <p>Uses predicate <code>in</code></p> <p>Uses predicate <code>adjacent</code> or <code>near</code> or <code>touch</code></p> <p>Uses object category <code>balls</code></p> <p>Uses object category <code>receptacles</code></p> |
| <pre>(preference pillowMovedToRoomCenter (exists (?o - pillow) (then (once (and (agent_holds ?o))) (hold (and (in_motion ?o) (not (agent_holds ? o)))) (once (and (not (in_motion ?o)) (near room_center ?o) (exists (?o1 ?o2 ?o3 - game_object) (and (same_color ?o1 pink) (near room_center ?o1) (same_color ?o2 blue) (near room_center ?o2) (same_color ?o3 brown) (near room_center ?o3))))))))) </pre> | <p>This preference is satisfied when:</p> <ul style="list-style-type: none"> -first, the agent picks up a pillow -next, the agent throws the pillow and it is no longer being held by the agent -finally, the pillow stops moving near the center of the room, and there are three other objects near the center of the room as well: one that is pink, one that is blue, and one that is brown | <p>Uses predicate <code>agent_holds</code> or <code>in_motion</code></p> <p>Uses predicate <code>adjacent</code> or <code>near</code> or <code>touch</code></p> <p>Uses object category <code>furniture</code> or <code>room_features</code></p> <p>Uses object category <code>small_objects</code> or <code>large_objects</code> or <code>any_object</code></p> |

Algorithm 1 Fitness model training loop

Input: Real games $\mathcal{D}^+ \in \mathbb{R}^{N \times \times 1F}$, regrown games $\mathcal{D}^- \in \mathbb{R}^{N \times M \times F}$
Input: Fitness model $f_\theta : \mathbb{R}^F \rightarrow \mathbb{R}$, optimizer
 N positive examples, M negatives generated per positive, B batch size, F features, K negatives sampled per positive in each epoch, P plateau epochs
Output: Converged fitness model W_θ
best model \leftarrow None
best loss $\leftarrow \infty$
last improvement epoch $\leftarrow -1$
for epoch i **do**
 \triangleright Assign negatives randomly to each positive
 Shuffle the first two dimensions of \mathcal{D}^-
 \triangleright Reorder the positives in each epoch
 Shuffle the first dimension of \mathcal{D}^+
 for each batch **do**
 $X^+ \leftarrow$ the next B positives $\triangleright X^+: B \times 1 \times F$
 $X^- \leftarrow K$ sampled negatives for each positive $\triangleright X^-: B \times K \times F$
 $X \leftarrow \text{concat}(X^+, X^-)$ $\triangleright X: B \times (1 + K) \times F$
 $Y \leftarrow f_\theta(X)$ $\triangleright Y: B \times (1 + K)$
 $L \leftarrow \text{softmax loss}(Y)$ $\triangleright L: \text{scalar}$
 Take backward step on loss and optimizer step
 end for
 epoch validation losses $\leftarrow []$
 for each batch in validation **do**
 <the above procedure without the optimizer steps>
 <append each batch’s loss to epoch validation losses>
 end for
 epoch loss $\leftarrow \text{mean}(\text{epoch validation losses})$
 if epoch loss $<$ best loss **then**
 best model \leftarrow copy of $f_\theta(X)$
 best loss \leftarrow epoch loss
 last improvement epoch $\leftarrow i$
 else if $i - \text{last improvement epoch} > P$ **then**
 break
 end if
end for
return best model

preferences (up to 4). In total, this set of behavioral characteristics allows for an archive size of 2000 games, of which 20 have one preference (matching one of the 9 exemplars or matching none of them, with and without a setup component), 110 have two preferences, 440 have three preferences, and 1430 have four preferences.

In addition, we add one more “pseudo behavioral characteristic” that explicitly captures a few general coherence properties of games – specifically features that we expect either *all* plausibly human-generated games to either exhibit or *none* of them to exhibit. While these features are also used by our learned fitness function, we use this behavioral characteristic as a sort of first-stage filter: if a game fails to meet these criteria, then it cannot reasonably be said to be “human-quality,” regardless of its fitness evaluation. For all reported games, we ensure that each of the criteria are satisfied. The criteria included in this behavioral characteristic include whether all variables are defined / used in preferences, whether all preferences are used in either terminal or scoring conditions, and whether the game avoids a set of grammatical but obviously nonsensical or redundant expressions. There are a total of 21 features used in this behavioral characteristic. This “pseudo behavioral characteristic” doubles the size of the archive (from 2000 games to 4000 games), though we never evaluate any game from the half of the archive in which this feature is false.

We begin the MAP-Elites algorithm by generating 1024 random games from the PCFG. We then sort each of the games in descending order of fitness and add them to the archive until either (a) every

possible value of each behavioral characteristic is represented by at least one game (note that this is not the same as every possible *combination* of behavioral characteristic values being represented), or **(b)** at least 128 cells of the archive are occupied.

We run MAP-Elites for 8192 “generations,” where each generation consists of 750 potential updates in which we randomly select a parent game, sample a mutation operator to apply, and potentially add the resulting mutated game to the archive.

DSL to natural language back-translation

In order to prepare games for human evaluation, we convert them from the DSL to natural language in a multi-stage process. In order to ensure consistency, we perform this back translation on both generated games and the real games (as opposed to using the original human-authored descriptions).

In the first stage of back-translation, a rule based system converts the DSL into templated language by concretely describing the definition of each predicate and grammatical rule. For instance, the expression `(once (and (agent_holds ?d) (adjacent ?p agent)))` is converted to “there is a state where (the agent is holding ?d) and (?p is adjacent to agent).” Each of the game’s setup conditions, preferences, terminal conditions, and scoring rules are rendered in this form, which also includes the mapping from variable names (e.g. ?d) to the types of objects that can occupy the variable (e.g. “dodgeball”). An example of a game’s preferences described in this form is presented below:

The preferences of the game are:

```

----Preference 1----
The variables required by this preference are:
-p of type pyramid_block
-d of type dodgeball
-h of type hexagonal_bin

This preference is satisfied when:
- first, there is a state where (the agent is holding ?d) and (?p is adjacent to agent)
- next, there is a sequence of one or more states where (it's not the case that the agent is holding ?d) and (?d is in motion)
- finally, there is a state where (it's not the case that ?d is in motion) and (?d is inside of ?h)

----Preference 2----
The variables required by this preference are:
-b of type building
-l of type cube_block
-f of type flat_block

This preference is satisfied when:
- in the final game state, (?f is used in the setup), (?f is inside of ?b), and (?l is inside of ?b)

```

Next, we use the GPT-4 large language model (LLM) [52] to simplify the templated description into a more naturalistic form (specifically `gpt-4-1106-preview`). The objective of this stage is to re-write any unclear formulations generated by the initial procedure and to replace abstract variable names with their actual referents. We convert each section of the game separately, using a similar prompt for each. The prompt begins with the following message:

“Your task is to convert a templated description of a game’s <setup / rules / terminal conditions / scoring conditions> into a natural language description. Do not change the content of the template, but you may rewrite and reorder the information in any way you think is necessary in order for a human to understand it. Use simple language and verbs that would be familiar to a human who has never played this game before.”

We then include 10 examples of this kind of translation taken from the set of human games not used in our experiments. An example of the same preferences in this simplified form is presented below:

```

The preferences of the game are:

----Preference 1----
This preference is satisfied when:
-first, the agent holds a dodgeball while standing next to a pyramid block
-next, the agent throws the dodgeball
-finally, the dodgeball lands inside a hexagonal bin and stops moving

----Preference 2----
This preference is satisfied when:
-at the end of the game, a flat block is used in the setup of a building and both a cube block and the flat block are inside the building

```

Finally, we use the LLM again to collect the separate descriptions of each section into one a single block, further simplifying the language and expressions. The prompt is similar to that used in the previous stage, and again is followed by 10 selected examples:

“Your task is to combine and simplify the description of a game’s rules. Do not change the content of the rules by either adding or removing information, but you may rewrite and reorder the information in any way you think is necessary in order for a human to understand it. Use simple language and verbs that would be familiar to a human who has never played this game before. DO describe preferences carefully, such that a player reading the description can easily play the game. DO NOT include explicit references to a game’s preferences (i.e. "Preference 1" or "Preference 2"). DO NOT include descriptions of setup or terminal conditions if they do not appear in the game.”

Examples of complete translations are available in Figure 4 and Figure 5.

E Human Evaluation Data Analysis

E.1 Detailed human evaluation results

Table SI-2: Human evaluation result summary

| Attribute | Mean score by category | | | Real vs. Matched | Real vs. Unmatched | Matched vs. Unmatched |
|----------------------------|------------------------|---------|-----------|----------------------|---------------------------|---------------------------|
| | Real | Matched | Unmatched | U-stat, P-value | U-stat, P-value | U-stat, P-value |
| <i>Understandable</i> | 3.943 | 3.923 | 3.331 | 45088.0, $P = 0.906$ | 55921.5, $P < 1e-5^{***}$ | 55846.0, $P < 1e-5^{***}$ |
| <i>Fun to play</i> | 2.522 | 2.430 | 2.068 | 46752.5, $P = 0.352$ | 54040.5, $P < 1e-5^{***}$ | 52539.5, $P < 1e-3^{***}$ |
| <i>Fun to watch</i> | 2.385 | 2.313 | 2.024 | 46169.0, $P = 0.519$ | 51636.5, $P < 1e-3^{***}$ | 50515.0, $P = 0.001^{**}$ |
| <i>Helpful[†]</i> | 2.997 | 2.987 | 2.840 | 44802.0, $P = 0.982$ | 47372.5, $P = 0.078$ | 47559.0, $P = 0.075$ |
| <i>Difficult</i> | 2.582 | 2.660 | 2.676 | 42921.5, $P = 0.326$ | 42218.5, $P = 0.419$ | 44081.0, $P = 0.947$ |
| <i>Creative</i> | 2.318 | 2.213 | 2.143 | 47036.0, $P = 0.282$ | 48286.0, $P = 0.025^*$ | 46615.0, $P = 0.182$ |
| <i>Human-like</i> | 2.813 | 2.670 | 2.119 | 47698.0, $P = 0.167$ | 58679.0, $P < 1e-5^{***}$ | 55434.5, $P < 1e-5^{***}$ |

Evaluators don’t distinguish between participant-created real and matched model games, but do distinguish unmatched games from both. Participants responded to seven Likert questions on a 5-point scale, one for each attribute in the first column (see Human evaluation methods for additional details). We report the Mann-Whitney U test [46] for differences in outcomes, appropriate for ordinal data. *: $P < 0.05$, **: $P < 0.01$, ***: $P < 0.001$

†: The full measure description is “Helpful for interacting with the simulated environment.”

E.2 Mixed-effect model analyses

In section , we briefly describe the mixed effect model we fit to analyze our human evaluation results and analyze the learned regression weights for the fitness function. Here, we build on this analysis to examine the extent to which accounting for the mediating effect of fitness scores, changes our previous observations regarding the differences between groups. Using the unmatched group as the baseline, the regression coefficients β_{matched} and β_{real} quantify these differences for each measure. We find statistically significant differences for the matched group (i.e. $\beta_{\text{matched}} > 0$) for ratings of fun to play, fun to watch, helpfulness, and human likeness. Similarly, we observe statistically significant differences ($\beta_{\text{real}} > 0$) for ratings of understandability, fun to play and watch, and human likeness. Finally, using the marginal (least-squares) means method[41], we directly compare the matched and real categories and again find no statistically significant differences (see Human evaluation methods for additional details and Supplementary Table SI-4 below for the full results).

Table SI-3: Mixed model result summary

| Attribute | Fitness | | | Matched | | | Real | | |
|-----------------------------|--------------------------|--------|------------------|--------------------------|-------|------------------|-----------------------|--------|------------------|
| | β_{fitness} | Z | P-value | β_{matched} | Z | P-value | β_{real} | Z | P-value |
| <i>Understandable</i> | 0.846 | 5.625 | $P < 1e-5^{***}$ | 0.525 | 1.766 | $P = 0.078$ | 1.151 | 4.036 | $P < 1e-3^{***}$ |
| <i>Fun to play</i> | 0.396 | 2.936 | $P = 0.003^{**}$ | 0.629 | 2.298 | $P = 0.022^*$ | 1.059 | 4.021 | $P < 1e-3^{***}$ |
| <i>Fun to watch</i> | 0.191 | 1.469 | $P = 0.142$ | 0.641 | 2.414 | $P = 0.016^*$ | 0.912 | 3.547 | $P < 1e-3^{***}$ |
| <i>Helpful</i> [†] | -0.189 | -2.163 | $P = 0.031^*$ | 0.349 | 2.048 | $P = 0.041^*$ | 0.232 | 1.441 | $P = 0.15$ |
| <i>Difficult</i> | -0.588 | -3.443 | $P < 1e-3^{***}$ | 0.363 | 1.029 | $P = 0.304$ | -0.250 | -0.740 | $P = 0.46$ |
| <i>Creative</i> | -0.486 | -3.191 | $P = 0.001^{**}$ | 0.551 | 1.776 | $P = 0.076$ | 0.438 | 1.467 | $P = 0.142$ |
| <i>Human-like</i> | 0.570 | 4.316 | $P < 1e-3^{***}$ | 0.837 | 3.128 | $P = 0.002^{**}$ | 1.446 | 5.597 | $P < 1e-5^{***}$ |

Fitness scores significantly predict several attributes, including understandability and human-likeness. Fitness scores show (statistically) significant positive effects on the understandability, fun to play, and human-likeness attributes, and significant negative effects on the difficulty and creativity questions. Accounting for the role of fitness, the `matched` group membership shows a significant effect only on human likeness. The `real` group shows significant effects on understandability, fun to play to watch, and human likeness. *: $P < 0.05$, **: $P < 0.01$, ***: $P < 0.001$

†: The full measure description is “Helpful for interacting with the simulated environment.”

E.3 Marginal Means Analysis

Table SI-4: Mixed model marginal means comparison summary

| Attribute | Real – Matched | | | Real – Unmatched | | | Matched – Unmatched | | |
|-------------------------------|----------------|--------|-------------|------------------|--------|------------------|---------------------|-------|------------------|
| | Diff | Z | P-value | Diff | Z | P-value | Diff | Z | P-value |
| <i>Understandable</i> ↑ | 0.626 | 2.055 | $P = 0.100$ | 1.151 | 4.036 | $P < 1e-3^{***}$ | 0.525 | 1.766 | $P = 0.181$ |
| <i>Fun to play</i> ↑ | 0.430 | 1.577 | $P = 0.256$ | 1.059 | 4.021 | $P < 1e-3^{***}$ | 0.629 | 2.298 | $P = 0.056$ |
| <i>Fun to watch</i> ↑ | 0.271 | 1.025 | $P = 0.561$ | 0.912 | 3.547 | $P = 0.001^{**}$ | 0.641 | 2.414 | $P = 0.042^*$ |
| <i>Helpful</i> [†] ↑ | -0.117 | -0.701 | $P = 0.763$ | 0.232 | 1.441 | $P = 0.32$ | 0.349 | 2.048 | $P = 0.101$ |
| <i>Difficult</i> ↓ | -0.613 | -1.725 | $P = 0.196$ | -0.250 | -0.740 | $P = 0.74$ | 0.363 | 1.029 | $P = 0.559$ |
| <i>Creative</i> ↑ | -0.113 | -0.364 | $P = 0.93$ | 0.438 | 1.467 | $P = 0.307$ | 0.551 | 1.776 | $P = 0.178$ |
| <i>Human-like</i> ↑ | 0.609 | 2.299 | $P = 0.056$ | 1.446 | 5.597 | $P < 1e-5^{***}$ | 0.837 | 3.128 | $P = 0.005^{**}$ |

We use the method of marginal (least-squares) means [41] to estimate the mean score for each attribute in each category, holding fitness constant. None of the comparisons between the `real` and `matched` groups are significant, and several (though not all) of the previously significant comparisons remain significant. *: $P < 0.05$, **: $P < 0.01$, ***: $P < 0.001$

†: The full measure description is “Helpful for interacting with the simulated environment.”

In most measures, higher scores are better, indicated by the ↑, other than *Difficult* ↓, in which 3 means “appropriately difficult”, and scores below and above indicate too easy and too hard respectively.

E.4 Matched-real game similarity analysis

To functionally measure similarity, we leverage the fact that our goal programs are interpretable and automatically evaluate them on all gameplay interactions generated by participants in our first experiment. For each `matched` game and its corresponding `real` counterpart, we measure the number of interactions that fulfill a gameplay element in only the `matched` game, only the `real` game, or both. While some pairs of games have their elements fulfilled by the same interactions (suggesting functional similarity), most pairs are not — under 25% (7/30) share more than half of their relevant interactions. Furthermore, the average Jaccard similarity between the sets of relevant interactions for the `matched` and `real` game is only 0.347 and the median similarity is 0.180 (identical games would score 1.0, entirely dissimilar games 0; and see summary in Figure ED-6 and methodological details in Sample similarity comparison methods).

F Model Ablations

F.1 Common Sense Ablation

The domain-specific language we use is underconstrained—many expressions that are grammatical either make no sense at all (e.g., checking a bin is in a ball, rather than a ball in a bin) or violate

intuitive physical common sense (e.g., creating a game stacking balls, as opposed to stacking blocks). We primarily operationalize the concept of physical common sense using two of our fitness features, discussed in Fitness function methods and Supplementary information B. Both use a dataset of interaction traces (see Dataset collection methods) to estimate the feasibility of predicate role-filler expressions, by computing the proportion of predicate expressions (and the object types they operate over) that have appeared at least once over the set of interactions of users with the environment. While this condition is not necessary (as it is unlikely experiment participants explored every feasible configuration of objects in the environment), it is sufficient to determine feasibility and, therefore, serves as a good proxy for intuitive common sense. The first feature operates over individual predicates, e.g. estimating that (`on desk ball`) is more likely than (`on desk bed`). The second feature operates on logical expressions over predicates, and might help catch contradictory predicates that are independently feasible, such as (`and (on desk ball) (on bed ball)`), that is, the ball might feasibly be on the desk or on the bed, but not on both.

We know that these features are helpful for our model, as the individual predicate version of these features has the third highest weight of all features that predict real human-generated games (see Supplementary information B.1 for details). To further evaluate the importance of these features, we fit a version of our fitness model that has no access to them, and use it as the objective for our model. Unsurprisingly, when we evaluate samples from this ablated model on the full fitness function (with the interaction trace features), they have statistically significantly lower fitness scores than the samples from the full model (matched-pairs t -test matching by archive cells, $t = -32.66$, $P < 1e-10$). To offer a more fair comparison, we use the full “reward machine” and dataset of play traces. We assign a binary score to each game from the full and ablated models, 1 if each game component (gameplay preferences and the setup section (if one exists)) is satisfied at least once over the dataset, either in the same trace or in different traces. If at least one game component is never satisfied, we assign a score of 0. We find that 1515 (75.75%) of the games in the full model score 1, while only 584 (29.20%) in the ablated model do. This difference is, as expected, also statistically significant (matched-pairs t -test, $t = -33.29$, $P < 1e-10$). We conclude that intuitive physical common sense is helpful to our model, as allowing our model to approximate the physical sensibility of predicates helps the model generate games with components that have been satisfied by our participants.

F.2 Compositionality Ablation

Evaluating the role of compositionality in our model is challenging as the model operates on a domain-specific language that is inherently highly compositional. Given the nature of program representations, it’s difficult to imagine a non-compositional counterfactual DSL to compare to — so we cannot compare to an entirely non-compositional model. Instead, we ablate by varying how compositional we allow our MAP-Elites mutation operators to be. The primary operator embodying compositionality is the crossover operator, which samples two programs from the MAP-Elites archive, randomly selects exchangeable sub-trees from both programs, and creates new candidates with these trees swapped between the programs. We also implemented several custom operators (beyond the evolutionary programming staples of mutation, insertion, deletion, and crossover). Many of these implement targeted variations of crossover that we considered to be plausible higher-level changes a person might make to a game they are creating, such as sampling a preference from another game and then changing the preference’s initial or terminal conditions. We report two ablations, one (“No Custom Ops”) where we omit the custom operators we implemented (keeping only regrowth, insertion, deletion, and crossover), and a second (“No Custom Ops, No Crossover”) where we also remove the crossover operation. We keep all other model details identical, crucially both the set of behavioral characteristics and fitness function, allowing us to directly compare the fitness values of games in the archives in the ablated models.

We visualize the results of these in Figure SI-2. While removing our custom operators appears to slightly *increase* the mean fitness of exemplars (Figure SI-2a, orange), removing the crossover operation drastically decreases the fitness of games in that model’s archive (Figure SI-2a, green). This provides evidence that allowing our search procedure to take advantage of the compositionality in our domain is greatly beneficial in generating high-quality samples across our archive. If our custom operators do not increase mean fitness, what impact do they have? To quantify this question, we evaluated samples from the ablated models through the full “reward machine.” For each sample, we counted how many of the participant interaction traces saw the participant fulfilling one or more gameplay elements from the sample. In other words, how many participants (unwittingly or otherwise)

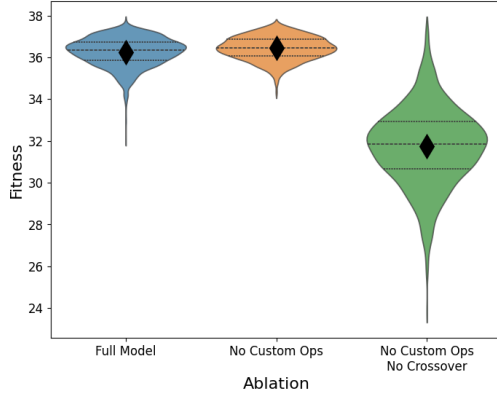
fulfilled at least part of the model-generated goal program? We find that the custom operators help increase this number – samples generated from our full model show the highest number of relevant traces (Figure SI-2b, blue). This could have two interpretations: one is that the custom operators push more goals toward higher feasibility. Another is that this behavior is a form of mode-seeking that helps the model generate goal programs that capture more common behaviors, as opposed to more meaningful variability. In all, we take this as an effect that crossover is crucial to generating fit samples across our MAP-Elites archive, with some cost to diversity which our custom operators help reduce.

F.3 Coherence ablation

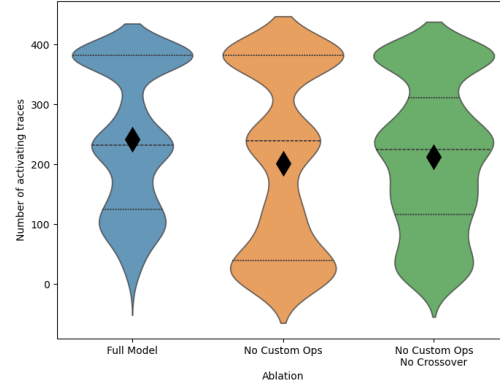
As we iterated on earlier versions of our model, we discovered some ‘softer’, higher-level issues repeatedly surfacing in model-generated goal programs. Even after implementing features that helped the model avoid some types of low-level mistakes (such as instantiating variables or preferences and never referencing them), and introducing approximations to intuitive physical common sense (discussed above), some aspects of the generated games remained incoherent. A lower-level example might be disjointedness in the arguments of temporal modals. Consider, for instance, a preference whose modals translate to natural language as “start with a state where the agent holds a ball, then find a collection of states where a block is on the bed, and finish with a state where the bin is upside down”. The awkwardness in explaining this perfectly grammatical preference (program below) is that each modal ((`once` . . .), (`hold` . . .)) refers to a distinct set of objects, and so it feels unnatural to specify a sequential temporal preference over them.

```
(preference preference0
  (exists (?v0 – hexagonal_bin ?v1 – ball ?v2 – block)
    (then
      (once (agent_holds ?v1))
      (hold (on desk ?v2))
      (once (object_orientation ?v0 upside_down))
    )))
```

We observed similar, higher-level issues regarding coherence between different gameplay preferences (do they use the same objects and predicates, or distinct sets?). Specifically, we observed cases where game scoring conditions and ending conditions have nothing to do with each other. A game might specify that it ends after a ball has been thrown five times, with points scored for every block placed on the desk. There is nothing wrong per se with this specification, but it feels unnatural—we would expect either the ball-throwing to contribute to scoring, or the block-stacking to allow the game to end, or both. We wrote a collection of fitness features to try to capture occurrences of such incoherence (see `game_element_disjointness` in Supplementary information B). We have some indication that these features are important from observing that our fitness model assigns one of them the third-largest negative weight (predictive of corrupted, negative games). To ablate the effect of this feature group, we perform an ablation similarly to the common sense ablation reported above—we fit a fitness model without these features and use it to guide our MAP-Elites search. As a first sanity check, we compute fitness scores under the full fitness model for games generated by the ablated model. We find that scores in the ablated model are consistently lower (matched-pairs *t*-test, $t = -26.99$, $P < 1e-10$), indicating that without access to these features, our model would generate programs that violate these coherence considerations. We also evaluate games from this ablated model using the ‘reward machine’ and play traces dataset, as we did above. As before, 1515 (75.75%) of games in the full model have every component satisfied, while only 1224 (61.2%) in the ablated model do. This difference is also statistically significant (matched-pairs *t*-test, $t = -9.73$, $P < 1e-10$).



(a) Removing crossover drastically lowers fitness values. We plot, for each game generated by a model, its fitness score under the full fitness function. **Left:** The distribution of fitness scores in our full model. **Middle:** Removing the custom operators has little effect on the distribution of fitness scores. **Right:** Removing crossover drastically lowers the fitness scores of model samples.



(b) Removing custom operators lowers mean trace coverage; removing crossover undoes some of the effect. We measure, for each game generated by a model, how many participant interaction traces fulfill at least one gameplay element. **Left:** Of the ablations reported, our full model shows the highest number of active traces. **Middle:** removing our custom mutation operators lowers the mean number of active traces. **Right:** Removing crossover as well undoes some of the effect of removing the custom operators.

Figure SI-2: The crossover operator helps generate fit goals, while the custom operators help generate solutions with higher trace coverage. **Left:** removing the custom operators does hurt mean fitness scores; removing the crossover operator does. **Right:** removing the custom operators leads the model to generate samples covering fewer participant interaction traces on average. This could be evidence of lower feasibility (more samples in the “no custom ops” model are active in barely a few traces) or of mode seeking (more samples in the full model are active in a very high number of traces).

G Full Domain Specific Language Description

G.1 DSL Grammar Definitions

A game is defined by a name, and is expected to be valid in a particular domain, also referenced by a name. A game is defined by four elements, two of them mandatory, and two optional. The mandatory ones are the $\langle constraints \rangle$ section, which defines gameplay preferences, and the $\langle scoring \rangle$ section, which defines how gameplay preferences are counted to arrive at a score for the player in the game. The optional ones are the $\langle setup \rangle$ section, which defines how the environment must be prepared before gameplay can begin, and the $\langle terminal \rangle$ conditions, which specify when and how the game ends.

```

<game> ::= (define (game <ID>)
  (:domain <ID>)
  (:setup <setup>)
  (:constraints <constraints>)
  (:terminal <terminal>)
  (:scoring <scoring>)
)

```

$\langle id \rangle ::= /[a-z0-9][a-z0-9]+/$ # a letter or digit, followed by one or more letters, digits, or dashes

We will now proceed to introduce and define the syntax for each of these sections, followed by the non-grammar elements of our domain: predicates, functions, and types. Finally, we provide a mapping between some aspects of our gameplay preference specification and linear temporal logic (LTL) operators.

G.1.1 Setup

The setup section specifies how the environment must be transformed from its deterministic initial conditions to a state gameplay can begin at. Currently, a particular environment room always appears in the same initial conditions, in terms of which objects exist and where they are placed. Participants in our experiment could, but did not have to, specify how the room must be setup so that their game could be played.

The initial $\langle setup \rangle$ element can expand to conjunctions, disjunctions, negations, or quantifications of itself, and then to the $\langle setup-statement \rangle$ rule. $\langle setup-statement \rangle$ elements specify two different types of setup conditions: either those that must be conserved through gameplay ('game-conserved'), or those that are optional through gameplay ('game-optional'). These different conditions arise as some setup elements must be maintain through gameplay (for example, a participant specified to place a bin on the bed to throw balls into, it shouldn't move unless specified otherwise), while other setup elements can or must change (if a participant specified to set the balls on the desk to throw them, an agent will have to pick them up (and off the desk) in order to throw them).

Inside the $\langle setup-statement \rangle$ tags we find $\langle super-predicate \rangle$ elements, which are logical operations and quantifications over other $\langle super-predicate \rangle$ elements, function comparisons ($\langle function-comparison \rangle$), which like predicates also resolve to a truth value), and predicates ($\langle predicate \rangle$). Function comparisons usually consist of a comparison operator and two arguments, which can either be the evaluation of a function or a number. The one exception is the case where the comparison operator is the equality operator (=), in which case any number of arguments can be provided. Finally, the $\langle predicate \rangle$ element expands to a predicate acting on one or more objects or variables. For a full list of the predicates we found ourselves using so far, see Appendix G.2.1.

$\langle setup \rangle ::= (\text{and } \langle setup \rangle \langle setup \rangle^+) \#$ A setup can be expanded to a conjunction, a disjunction, a quantification, or a setup statement (see below).

- | (or $\langle setup \rangle \langle setup \rangle^+$)
- | (not $\langle setup \rangle$)
- | (exists ($\langle typed\ list(variable) \rangle$)) $\langle setup \rangle$)
- | (forall ($\langle typed\ list(variable) \rangle$)) $\langle setup \rangle$)
- | $\langle setup-statement \rangle$

$\langle setup-statement \rangle ::= \#$ A setup statement specifies that a predicate is either optional during gameplay or must be preserved during gameplay.

- | (game-conserved $\langle super-predicate \rangle$)
- | (game-optional $\langle super-predicate \rangle$)

$\langle super-predicate \rangle ::= \#$ A super-predicate is a conjunction, disjunction, negation, or quantification over another super-predicate. It can also be directly a function comparison or a predicate.

- | (and $\langle super-predicate \rangle^+$)
- | (or $\langle super-predicate \rangle^+$)
- | (not $\langle super-predicate \rangle$)
- | (exists ($\langle typed\ list(variable) \rangle$)) $\langle super-predicate \rangle$)
- | (forall ($\langle typed\ list(variable) \rangle$)) $\langle super-predicate \rangle$)
- | $\langle f-comp \rangle$
- | $\langle predicate \rangle$

$\langle function-comparison \rangle ::= \#$ A function comparison: either comparing two function evaluations, or checking that two ore more functions evaluate to the same result.

- | ($\langle comp-op \rangle \langle function-eval-or-number \rangle \langle function-eval-or-number \rangle$)
- | (= $\langle function-eval-or-number \rangle^+$)

$\langle comp-op \rangle ::= \langle | \langle = | = | \rangle \rangle = \#$ Any of the comparison operators.

$\langle function-eval-or-number \rangle ::= \langle function-eval \rangle | \langle comparison-arg-number \rangle$

$\langle comparison-arg-number \rangle ::= \langle number \rangle$

$\langle number \rangle ::= /-?\d*\.?\d+/\ #$ A number, either an integer or a float.

⟨function-eval⟩ ::= # See valid expansions in a separate section below

⟨variable-list⟩ ::= (*⟨variable-def⟩*)⁺ # One or more variables definitions, enclosed by parentheses.

⟨variable-def⟩ ::= *⟨variable-type-def⟩* | *⟨color-variable-type-def⟩* | *⟨orientation-variable-type-def⟩* | *⟨side-variable-type-def⟩* # Colors, sides, and orientations are special types as they are not interchangeable with objects.

⟨variable-type-def⟩ ::= *⟨variable⟩*⁺ - *⟨type-def⟩* # Each variable is defined by a variable (see next) and a type (see after).

⟨color-variable-type-def⟩ ::= *⟨color-variable⟩*⁺ - *⟨color-type-def⟩* # A color variable is defined by a variable (see below) and a color type.

⟨orientation-variable-type-def⟩ ::= *⟨orientation-variable⟩*⁺ - *⟨orientation-type-def⟩* # An orientation variable is defined by a variable (see below) and an orientation type.

⟨side-variable-type-def⟩ ::= *⟨side-variable⟩*⁺ - *⟨side-type-def⟩* # A side variable is defined by a variable (see below) and a side type.

⟨variable⟩ ::= $\Lambda?[a-w][a-z0-9]^*$ # a question mark followed by a lowercase a-w, optionally followed by additional letters or numbers.

⟨color-variable⟩ ::= $\Lambda?x[0-9]^*$ # a question mark followed by an x and an optional number.

⟨orientation-variable⟩ ::= $\Lambda?y[0-9]^*$ # a question mark followed by an y and an optional number.

⟨side-variable⟩ ::= $\Lambda?z[0-9]^*$ # a question mark followed by an z and an optional number.

⟨type-def⟩ ::= *⟨object-type⟩* | *⟨either-types⟩* # A variable type can either be a single name, or a list of type names, as specified below

⟨color-type-def⟩ ::= *⟨color-type⟩* | *⟨either-color-types⟩* # A color variable type can either be a single color name, or a list of color names, as specified below

⟨orientation-type-def⟩ ::= *⟨orientation-type⟩* | *⟨either-orientation-types⟩* # An orientation variable type can either be a single orientation name, or a list of orientation names, as specified below

⟨side-type-def⟩ ::= *⟨side-type⟩* | *⟨either-side-types⟩* # A side variable type can either be a single side name, or a list of side names, as specified below

⟨either-types⟩ ::= (either *⟨object-type⟩*)⁺

⟨either-color-types⟩ ::= (either *⟨color⟩*)⁺

⟨either-orientation-types⟩ ::= (either *⟨orientation⟩*)⁺

⟨either-side-types⟩ ::= (either *⟨side⟩*)⁺

⟨object-type⟩ ::= *⟨name⟩*

⟨name⟩ ::= $/[A-Za-z][A-Za-z0-9_]+/$ # a letter, followed by one or more letters, digits, or underscores

⟨color-type⟩ ::= 'color'

⟨color⟩ ::= 'blue' | 'brown' | 'gray' | 'green' | 'orange' | 'pink' | 'purple' | 'red' | 'tan' | 'white' | 'yellow'

⟨orientation-type⟩ ::= 'orientation'

⟨orientation⟩ ::= 'diagonal' | 'sideways' | 'upright' | 'upside_down'

$\langle side-type \rangle ::= 'side'$

$\langle side \rangle ::= 'back' \mid 'front' \mid 'left' \mid 'right'$

$\langle predicate \rangle ::= \#$ See valid expansions in a separate section below

$\langle predicate-or-function-term \rangle ::= \langle object-name \rangle \mid \langle variable \rangle$ # A predicate or function term can either be an object name (from a small list allowed to be directly referred to) or a variable.

$\langle predicate-or-function-color-term \rangle ::= \langle color \rangle \mid \langle color-variable \rangle$

$\langle predicate-or-function-orientation-term \rangle ::= \langle orientation \rangle \mid \langle orientation-variable \rangle$

$\langle predicate-or-function-side-term \rangle ::= \langle side \rangle \mid \langle side-variable \rangle$

$\langle predicate-or-function-type-term \rangle ::= \langle object-type \rangle \mid \langle variable \rangle$

$\langle object-name \rangle ::= 'agent' \mid 'bed' \mid 'desk' \mid 'door' \mid 'floor' \mid 'main_light_switch' \mid 'mirror' \mid 'room_center' \mid 'rug' \mid 'side_table' \mid 'bottom_drawer' \mid 'bottom_shelf' \mid 'east_sliding_door' \mid 'east_wall' \mid 'north_wall' \mid 'south_wall' \mid 'top_drawer' \mid 'top_shelf' \mid 'west_sliding_door' \mid 'west_wall'$

G.1.2 Gameplay Preferences

The gameplay preferences specify the core of a game's semantics, capturing how a game should be played by specifying temporal constraints over predicates. The name for the overall element, $\langle constraints \rangle$, is inherited from the PDDL element with the same name.

The $\langle constraints \rangle$ elements expand into one or more preference definitions, which are defined using the $\langle pref-def \rangle$ element. A $\langle pref-def \rangle$ either expands to a single preference ($\langle preference \rangle$), or to a $\langle pref-forall \rangle$ element, which specifies variants of the same preference for different objects, which can be treated differently in the scoring section. A $\langle preference \rangle$ is defined by a name and a $\langle preference-quantifier \rangle$, which expands to an optional quantification (exists, forall, or neither), inside of which we find the $\langle preference-body \rangle$.

A $\langle preference-body \rangle$ expands into one of two options: The first is a set of conditions that should be true at the end of gameplay, using the $\langle at-end \rangle$ operator. Inside an $\langle at-end \rangle$ we find a $\langle super-predicate \rangle$, which like in the setup section, expands to logical operations or quantifications over other $\langle super-predicate \rangle$ elements, function comparisons, or predicates.

The second option is specified using the $\langle then \rangle$ syntax, which defines a series of temporal conditions that should hold over a sequence of states. Under a $\langle then \rangle$ operator, we find two or more sequence functions ($\langle seq-func \rangle$), which define the specific conditions that must hold and how many states we expect them to hold for. We assume that there are no unaccounted states between the states accounted for by the different operators – in other words, the $\langle then \rangle$ operators expects to find a sequence of contiguous states that satisfy the different sequence functions. The operators under a $\langle then \rangle$ operator map onto linear temporal logic (LTL) operators, see Appendix G.3 for the mapping and examples.

The $\langle once \rangle$ operator specifies a predicate that must hold for a single world state. If a $\langle once \rangle$ operators appears as the first operator of a $\langle then \rangle$ definition, and a sequence of states S_a, S_{a+1}, \dots, S_b satisfy the $\langle then \rangle$ operator, it could be the case that the predicate is satisfied before this sequence of states (e.g. by S_{a-1}, S_{a-2} , and so forth). However, only the final such state, S_a , is required for the preference to be satisfied. The same could be true at the end of the sequence: if a $\langle then \rangle$ operator ends with a $\langle once \rangle$ term, there could be other states after the final state (S_{b+1}, S_{b+2} , etc.) that satisfy the predicate in the $\langle once \rangle$ operator, but only one is required. The $\langle once-measure \rangle$ operator is a slight variation of the $\langle once \rangle$ operator, which in addition to a predicate, takes in a function evaluation, and measures the value of the function evaluated at the state that satisfies the preference. This function value can then be used in the scoring definition, see Appendix G.1.4.

A second type of operator that exists is the $\langle hold \rangle$ operator. It specifies that a predicate must hold true in every state between the one in which the previous operator is satisfied, and until one in which the next operator is satisfied. If a $\langle hold \rangle$ operator appears at the beginning or an end of a $\langle then \rangle$ sequence,

it can be satisfied by a single state, Otherwise, it must be satisfied until the next operator is satisfied. For example, in the minimal definition below:

```
(then
  (once (pred\_a))
  (hold (pred\_b))
  (once (pred\_c))
)
```

To find a sequence of states S_a, S_{a+1}, \dots, S_b that satisfy this $\langle then \rangle$ operator, the following conditions must hold true: (1) $pred_a$ is true at state S_a , (2) $pred_b$ is true in all states $S_{a+1}, S_{a+2}, \dots, S_{b-2}, S_{b-1}$, and (3) $pred_c$ is true in state S_b . There is no minimal number of states that the hold predicate must hold for.

The last operator is $\langle hold\text{-}while \rangle$, which offers a variation of the $\langle hold \rangle$ operator. A $\langle hold\text{-}while \rangle$ receives at least two predicates. The first acts the same as the predicate in a $\langle hold \rangle$ operator. The second (and third, and any subsequent ones), must hold true for at least one state while the first predicate holds, and must occur in the order specified. In the example above, if we substitute $(hold (pred_b))$ for $(hold\text{-}while (pred_b) (pred_d) (pred_e))$, we now expect that in addition to $pred_b$ being true in all states $S_{a+1}, S_{a+2}, \dots, S_{b-2}, S_{b-1}$, that there is some state $S_d, d \in [a + 1, b - 1]$ where $pred_d$ holds, and another state, $S_e, e \in [d + 1, b - 1]$ where $pred_e$ holds.

$\langle constraints \rangle ::= \langle pref\text{-}def \rangle \mid (\text{and } \langle pref\text{-}def \rangle^+)$ # One or more preferences.

$\langle pref\text{-}def \rangle ::= \langle pref\text{-}forall \rangle \mid \langle preference \rangle$ # A preference definitions expands to either a forall quantification (see below) or to a preference.

$\langle pref\text{-}forall \rangle ::= (\text{forall } \langle variable\text{-}list \rangle \langle preference \rangle)$ # this syntax is used to specify variants of the same preference for different objects, which differ in their scoring. These are specified using the $\langle pref\text{-}name\text{-}and\text{-}types \rangle$ syntax element's optional types, see scoring below.

$\langle preference \rangle ::= (\text{preference } \langle name \rangle \langle preference\text{-}quantifier \rangle)$ # A preference is defined by a name and a quantifier that includes the preference body.

$\langle preference\text{-}quantifier \rangle ::=$ # A preference can quantify existentially or universally over one or more variables, or none.

```
| (exists ((variable-list)) <preference-body>)
| (forall ((variable-list)) <preference-body>)
| <preference-body>
```

$\langle preference\text{-}body \rangle ::= \langle then \rangle \mid \langle at\text{-}end \rangle$

$\langle at\text{-}end \rangle ::= (\text{at-end } \langle super\text{-}predicate \rangle)$ # Specifies a prediicate that should hold in the terminal state.

$\langle then \rangle ::= (\text{then } \langle seq\text{-}func \rangle \langle seq\text{-}func \rangle^+)$ # Specifies a series of conditions that should hold over a sequence of states – see below for the specific operators ($\langle seq\text{-}func \rangle$ s), and Section 2 for translation of these definitions to linear temporal logic (LTL).

$\langle seq\text{-}func \rangle ::= \langle once \rangle \mid \langle once\text{-}measure \rangle \mid \langle hold \rangle \mid \langle hold\text{-}while \rangle$ # Four of these temporal sequence functions currently exist:

$\langle once \rangle ::= (\text{once } \langle super\text{-}predicate \rangle)$ # The predicate specified must hold for a single world state.

$\langle once\text{-}measure \rangle ::= (\text{once } \langle super\text{-}predicate \rangle \langle function\text{-}eval \rangle)$ # The predicate specified must hold for a single world state, and record the value of the function evaluation, to be used in scoring.

$\langle hold \rangle ::= (\text{hold } \langle super\text{-}predicate \rangle)$ # The predicate specified must hold for every state between the previous temporal operator and the next one.

$\langle hold\text{-}while \rangle ::= (\text{hold-while } \langle super\text{-}predicate \rangle \langle super\text{-}predicate \rangle^+)$ # The first predicate specified must hold for every state between the previous temporal operator and the next one. While it does, at least one state must satisfy each of the predicates specified in the second argument onward

For the full specification of the $\langle super\text{-}predicate \rangle$ element, see Appendix G.1.1 above.

G.1.3 Terminal Conditions

Specifying explicit terminal conditions is optional, and while some of our participants chose to do so, many did not. Conditions explicitly specified in this section terminate the game. If none are specified, a game is assumed to terminate whenever the player chooses to end the game.

The terminal conditions expand from the $\langle terminal \rangle$ element, which can expand to logical conditions on nested $\langle terminal \rangle$ elements, or to a terminal comparison. The terminal comparison ($\langle terminal-comp \rangle$) expands to one of three different types of comparisons: $\langle terminal-time-comp \rangle$, a comparison between the total time spent in the game ($\langle total-time \rangle$) and a time number token, $\langle terminal-score-comp \rangle$, a comparison between the total score ($\langle total-score \rangle$) and a score number token, or $\langle terminal-pref-count-comp \rangle$, a comparison between a scoring expression ($\langle scoring-expr \rangle$, see below) and a preference count number token. In most cases, the scoring expression is a preference counting operation.

$\langle terminal \rangle ::= \#$ The terminal condition is specified by a conjunction, disjunction, negation, or comparison (see below).
| (and $\langle terminal \rangle^+$)
| (or $\langle terminal \rangle^+$)
| (not $\langle terminal \rangle$)
| $\langle terminal-comp \rangle$

$\langle terminal-comp \rangle ::= \#$ We support three types of terminal comparisons:
| $\langle terminal-time-comp \rangle$
| $\langle terminal-score-comp \rangle$
| $\langle terminal-pref-count-comp \rangle$

$\langle terminal-time-comp \rangle ::= (\langle comp-op \rangle \langle total-time \rangle \langle time-number \rangle) \#$ The total time of the game must satisfy the comparison.

$\langle terminal-score-comp \rangle ::= (\langle comp-op \rangle \langle total-score \rangle \langle score-number \rangle) \#$ The total score of the game must satisfy the comparison.

$\langle terminal-pref-count-comp \rangle ::= (\langle comp-op \rangle \langle scoring-expr \rangle \langle preference-count-number \rangle) \#$ The number of times the preference specified by the name and types must satisfy the comparison.

$\langle time-number \rangle ::= \langle number \rangle \#$ Separate type so the we can learn a separate distribution over times than, say, scores.

$\langle score-number \rangle ::= \langle number \rangle$

$\langle preference-count-number \rangle ::= \langle number \rangle$

$\langle comp-op \rangle ::= \langle | \langle = | = | \rangle =$

For the full specification of the $\langle scoring-expr \rangle$ element, see Appendix G.1.4 below.

G.1.4 Scoring

Scoring rules specify how to count preferences (count once, once for each unique objects that fulfill the preference, each time a preference is satisfied, etc.), and the arithmetic to combine preference counts to a final score in the game.

A $\langle scoring-expr \rangle$ can be defined by arithmetic operations on other scoring expressions, references to the total time or total score (for instance, to provide a bonus if a certain score is reached), comparisons between scoring expressions ($\langle scoring-comp \rangle$), or by preference evaluation rules. Various preference evaluation modes can expand the $\langle preference-eval \rangle$ rule, see the full list and descriptions below.

$\langle scoring \rangle ::= \langle scoring-expr \rangle \#$ The scoring conditions maximize a scoring expression.

$\langle scoring-expr \rangle ::= \#$ A scoring expression can be an arithmetic operation over other scoring expressions, a reference to the total time or score, a comparison, or a preference scoring evaluation.

- | $\langle \text{scoring-external-maximize} \rangle$
- | $\langle \text{scoring-external-minimize} \rangle$
- | $\langle (\text{multi-op}) \langle \text{scoring-expr} \rangle^+ \rangle$ # Either addition or multiplication.
- | $\langle (\text{binary-op}) \langle \text{scoring-expr} \rangle \langle \text{scoring-expr} \rangle \rangle$ # Either division or subtraction.
- | $\langle (- \langle \text{scoring-expr} \rangle) \rangle$
- | $\langle \text{total-time} \rangle$
- | $\langle \text{total-score} \rangle$
- | $\langle \text{scoring-comp} \rangle$
- | $\langle \text{preference-eval} \rangle$
- | $\langle \text{scoring-number-value} \rangle$

$\langle \text{scoring-external-maximize} \rangle ::= (\text{external-forall-maximize } \langle \text{scoring-expr} \rangle)$ # For any preferences under this expression inside a (forall ...), score only for the single externally-quantified object that maximizes this scoring expression.

$\langle \text{scoring-external-minimize} \rangle ::= (\text{external-forall-minimize } \langle \text{scoring-expr} \rangle)$ # For any preferences under this expression inside a (forall ...), score only for the single externally-quantified object that minimizes this scoring expression.

$\langle \text{scoring-comp} \rangle ::=$ # A scoring comparison: either comparing two expressions, or checking that two ore more expressions are equal.

- | $\langle (\text{comp-op}) \langle \text{scoring-expr} \rangle \langle \text{scoring-expr} \rangle \rangle$
- | $\langle (= \langle \text{scoring-expr} \rangle^+) \rangle$

$\langle \text{preference-eval} \rangle ::=$ # A preference evaluation applies one of the scoring operators (see below) to a particular preference referenced by name (with optional types).

- | $\langle \text{count} \rangle$
- | $\langle \text{count-overlapping} \rangle$
- | $\langle \text{count-once} \rangle$
- | $\langle \text{count-once-per-objects} \rangle$
- | $\langle \text{count-measure} \rangle$
- | $\langle \text{count-unique-positions} \rangle$
- | $\langle \text{count-same-positions} \rangle$
- | $\langle \text{count-once-per-external-objects} \rangle$

$\langle \text{count} \rangle ::= (\text{count } \langle \text{pref-name-and-types} \rangle)$ # Count how many times the preference is satisfied by non-overlapping sequences of states.

$\langle \text{count-overlapping} \rangle ::= (\text{count-overlapping } \langle \text{pref-name-and-types} \rangle)$ # Count how many times the preference is satisfied by overlapping sequences of states.

$\langle \text{count-once} \rangle ::= (\text{count-once } \langle \text{pref-name-and-types} \rangle)$ # Count whether or not this preference was satisfied at all.

$\langle \text{count-once-per-objects} \rangle ::= (\text{count-once-per-objects } \langle \text{pref-name-and-types} \rangle)$ # Count once for each unique combination of objects quantified in the preference that satisfy it.

$\langle \text{count-measure} \rangle ::= (\text{count-measure } \langle \text{pref-name-and-types} \rangle)$ # Can only be used in preferences including a $\langle \text{once-measure} \rangle$ modal, maps each preference satisfaction to the value of the function evaluation in the $\langle \text{once-measure} \rangle$.

$\langle \text{count-unique-positions} \rangle ::= (\text{count-unique-positions } \langle \text{pref-name-and-types} \rangle)$ # Count how many times the preference was satisfied with quantified objects that remain stationary within each preference satisfaction, and have different positions between different satisfactions.

$\langle \text{count-same-positions} \rangle ::= (\text{count-same-positions } \langle \text{pref-name-and-types} \rangle)$ # Count how many times the preference was satisfied with quantified objects that remain stationary within each preference satisfaction, and have (approximately) the same position between different satisfactions.

$\langle \text{count-once-per-external-objects} \rangle ::= (\text{count-once-per-external-objects } \langle \text{pref-name-and-types} \rangle) \#$
 Similarly to count-once-per-objects, but counting only for each unique object or combination of objects quantified in the (forall ...) block including this preference.

$\langle \text{pref-name-and-types} \rangle ::= \langle \text{name} \rangle \langle \text{pref-object-type} \rangle^* \#$ The optional $\langle \text{pref-object-type} \rangle$ s are used to specify a particular instance of the preference for a given object, see the $\langle \text{pref-forall} \rangle$ syntax above.

$\langle \text{pref-object-type} \rangle ::= : \langle \text{type-name} \rangle \#$ The optional type name specification for the above syntax. For example, pref-name:dodgeball would refer to the preference where the first quantified object is a dodgeball.

$\langle \text{scoring-number-value} \rangle ::= \langle \text{number} \rangle$

G.2 Non-Grammar Definitions

G.2.1 Predicates

The following section describes the predicates we define. Predicates operate over a specified number of arguments, which can be variables or object names, and return a boolean value (true/false).

```
(above <arg1> <arg2>) [5 references] ; Is the first object above the second object?
(adjacent <arg1> <arg2>) [78 references] ; Are the two objects adjacent? [will probably be implemented as distance below
some threshold]
(adjacent_side <3 or 4 arguments>) [15 references] ; Are the two objects adjacent on the sides specified? Specifying a side
for the second object is optional, allowing to specify <obj1> <side1> <obj2> or <obj1> <side1> <obj2> <side2>
(agent_crouches ) [2 references] ; Is the agent crouching?
(agent_holds <arg1>) [327 references] ; Is the agent holding the object?
(between <arg1> <arg2> <arg3>) [7 references] ; Is the second object between the first object and the third object?
(broken <arg1>) [2 references] ; Is the object broken?
(equal_x_position <arg1> <arg2>) [2 references] ; Are these two objects (approximately) in the same x position? (in our
environment, x, z are spatial coordinates, y is the height)
(equal_z_position <arg1> <arg2>) [5 references] ; Are these two objects (approximately) in the same z position? (in our
environment, x, z are spatial coordinates, y is the height)
(faces <arg1> <arg2>) [6 references] ; Is the front of the first object facing the front of the second object?
(game_over ) [4 references] ; Is this the last state of gameplay?
(game_start ) [3 references] ; Is this the first state of gameplay?
(in <arg1> <arg2>) [121 references] ; Is the second argument inside the first argument? [a containment check of some sort,
for balls in bins, for example]
(in_motion <arg1>) [315 references] ; Is the object in motion?
(is_setup_object <arg1>) [13 references] ; Is this the object of the same type referenced in the setup?
(object_orientation <arg1> <arg2>) [14 references] ; Is the first argument, an object, in the orientation specified by the
second argument? Used to check if an object is upright or upside down
(on <arg1> <arg2>) [168 references] ; Is the second object on the first one?
(open <arg1>) [3 references] ; Is the object open? Only valid for objects that can be opened, such as drawers.
(opposite <arg1> <arg2>) [4 references] ; So far used only with walls, or sides of the room, to specify two walls opposite
each other in conjunction with other predicates involving these walls
(rug_color_under <arg1> <arg2>) [11 references] ; Is the color of the rug under the object (first argument) the color
specified by the second argument?
(same_color <arg1> <arg2>) [23 references] ; If two objects, do they have the same color? If one is a color, does the object
have that color?
(same_object <arg1> <arg2>) [7 references] ; Are these two variables bound to the same object?
(same_type <arg1> <arg2>) [14 references] ; Are these two objects of the same type? Or if one is a direct reference to a
type, is this object of that type?
(toggled_on <arg1>) [4 references] ; Is this object toggled on?
(touch <arg1> <arg2>) [48 references] ; Are these two objects touching?
```

G.2.2 Functions

he following section describes the functions we define. Functions operate over a specified number of arguments, which can be variables or object names, and return a number.

```
(building_size <arg1>) [2 references] ; Takes in an argument of type building, and returns how many objects comprise the
building (as an integer).
(distance <arg1> <arg2>) [114 references] ; Takes in two arguments of type object, and returns the distance between the two
objects (as a floating point number).
(distance_side <arg1> <arg2> <arg3>) [6 references] ; Similarly to the adjacent_side predicate, but applied to distance.
Takes in three or four arguments, either <obj1> <side1> <obj2> or <obj1> <side1> <obj2> <side2>, and returns the
distance between the first object on the side specified to the second object (optionally to its specified side).
(x_position <arg1>) [4 references] ; Takes in an argument of type object, and returns the x position of the object (as a
floating point number).
```

G.2.3 Types

The types are currently not defined as part of the grammar, other than the small list of *⟨object-name⟩* tokens that can be directly referred to, and are marked with an asterisk below. The following enumerates all expansions of the various *⟨type⟩* rules:

```
game_object [33 references] ; Parent type of all objects
agent* [90 references] ; The agent
building [20 references] ; Not a real game object, but rather, a way to refer to structures the agent builds
----- (* \textbf{Blocks} *) -----
block [28 references] ; Parent type of all block types:
bridge_block [11 references]
bridge_block_green [0 references]
bridge_block_pink [0 references]
bridge_block_tan [0 references]
cube_block [38 references]
cube_block_blue [8 references]
cube_block_tan [1 reference]
cube_block_yellow [8 references]
cylindrical_block [11 references]
cylindrical_block_blue [0 references]
cylindrical_block_green [0 references]
cylindrical_block_tan [0 references]
flat_block [5 references]
flat_block_gray [0 references]
flat_block_tan [0 references]
flat_block_yellow [0 references]
pyramid_block [13 references]
pyramid_block_blue [3 references]
pyramid_block_red [2 references]
pyramid_block_yellow [2 references]
tall_cylindrical_block [7 references]
tall_cylindrical_block_green [0 references]
tall_cylindrical_block_tan [0 references]
tall_cylindrical_block_yellow [0 references]
tall_rectangular_block [0 references]
tall_rectangular_block_blue [0 references]
tall_rectangular_block_green [0 references]
tall_rectangular_block_tan [0 references]
triangle_block [3 references]
triangle_block_blue [0 references]
triangle_block_green [0 references]
triangle_block_tan [0 references]
----- (* \textbf{Balls} *) -----
ball [40 references] ; Parent type of all ball types:
beachball [23 references]
basketball [18 references]
dodgeball [108 references]
dodgeball_blue [6 references]
dodgeball_red [4 references]
dodgeball_pink [8 references]
golfball [25 references]
golfball_green [3 references]
golfball_white [0 references]
----- (* \textbf{Colors} *) -----
color [6 references] ; Likewise, not a real game object, mostly used to refer to the color of the rug under an object
blue [6 references]
brown [5 references]
gray [0 references]
green [8 references]
orange [3 references]
pink [19 references]
purple [4 references]
red [8 references]
tan [2 references]
white [1 reference]
yellow [14 references]
----- (* \textbf{Furniture} *) -----
bed* [51 references]
blinds [2 references] ; The blinds on the windows
desk* [40 references]
desktop [6 references]
main_light_switch* [3 references] ; The main light switch on the wall
side_table* [4 references] ; The side table/nightstand next to the bed
shelf_desk [2 references] ; The shelves under the desk
----- (* \textbf{Large moveable/interactable objects} *) -----
book [11 references]
chair [18 references]
laptop [7 references]
pillow [14 references]
teddy_bear [14 references]
----- (* \textbf{Orientations} *) -----
diagonal [1 reference]
sideways [2 references]
upright [10 references]
upside_down [1 reference]
----- (* \textbf{Ramps} *) -----
ramp [0 references] ; Parent type of all ramp types:
curved_wooden_ramp [17 references]
triangular_ramp [10 references]
```

```

triangular_ramp_green [1 reference]
triangular_ramp_tan [0 references]
----- (* \textbf{Receptacles} *) -----
doggie_bed [26 references]
hexagonal_bin [123 references]
drawer [5 references] ; Either drawer in the side table
bottom_drawer* [0 references] ; The bottom of the two drawers in the nightstand near the bed.
top_drawer* [6 references] ; The top of the two drawers in the nightstand near the bed.
----- (* \textbf{Room features} *) -----
door* [9 references] ; The door out of the room
floor* [26 references]
mirror* [0 references]
poster* [0 references]
room_center* [0 references]
rug* [37 references]
shelf [10 references]
bottom_shelf* [1 reference]
top_shelf* [5 references]
sliding_door [2 references] ; The sliding doors on the south wall (big windows)
east_sliding_door* [1 reference] ; The eastern of the two sliding doors (the one closer to the desk)
west_sliding_door* [0 references] ; The western of the two sliding doors (the one closer to the bed)
wall [17 references] ; Any of the walls in the room
east_wall* [0 references] ; The wall behind the desk
north_wall* [1 reference] ; The wall with the door to the room
south_wall* [2 references] ; The wall with the sliding doors
west_wall* [3 references] ; The wall the bed is aligned to
----- (* \textbf{Small objects} *) -----
alarm_clock [8 references]
cellphone [6 references]
cd [6 references]
credit_card [1 reference]
key_chain [5 references]
lamp [2 references]
mug [3 references]
pen [2 references]
pencil [2 references]
watch [2 references]
----- (* \textbf{Sides} *) -----
back [3 references]
front [9 references]
left [3 references]
right [2 references]

```

G.3 Modal Definitions in Linear Temporal Logic

G.3.1 Linear Temporal Logic definitions

We offer a mapping between the temporal sequence functions defined in Appendix G.1.2 and linear temporal logic (LTL) operators. As we were creating this DSL, we found that the syntax of the *⟨then⟩* operator felt more convenient than directly writing down LTL, but we hope the mapping helps reason about how we see our temporal operators functioning. LTL offers the following operators, using φ and ψ as the symbols (in our case, predicates). Assume the following formulas operate sequence of states S_0, S_1, \dots, S_n :

- **Next**, $X\psi$: at the next timestep, ψ will be true. If we are at timestep i , then $S_{i+1} \vdash \psi$
 - **Finally**, $F\psi$: at some future timestep, ψ will be true. If we are at timestep i , then $\exists j > i : S_j \vdash \psi$
 - **Globally**, $G\psi$: from this timestep on, ψ will be true. If we are at timestep i , then $\forall j : j \geq i : S_j \vdash \psi$
 - **Until**, $\psi U \varphi$: ψ will be true from the current timestep until a timestep at which φ is true. If we are at timestep i , then $\exists j > i : \forall k : i \leq k < j : S_k \vdash \psi$, and $S_j \vdash \varphi$.
 - **Strong release**, $\psi M \varphi$: the same as until, but demanding that both ψ and φ are true simultaneously: If we are at timestep i , then $\exists j > i : \forall k : i \leq k \leq j : S_k \vdash \psi$, and $S_j \vdash \varphi$.
- Aside:* there's also a **weak until**, $\psi W \varphi$, which allows for the case where the second is never true, in which case the first must hold for the rest of the sequence. Formally, if we are at timestep i , if $\exists j > i : \forall k : i \leq k < j : S_k \vdash \psi$, and $S_j \vdash \varphi$, and otherwise, $\forall k \geq i : S_k \vdash \psi$. Similarly there's **release**, which is the similar variant of strong release. We're leaving those two as an aside since we don't know we'll need them.

G.3.2 Satisfying a *⟨then⟩* operator

Formally, to satisfy a preference using a *⟨then⟩* operator, we're looking to find a sub-sequence of S_0, S_1, \dots, S_n that satisfies the formula we translate to. We translate a *⟨then⟩* operator by translating

the constituent sequence-functions (*once*), (*hold*), (*while-hold*)¹ to LTL. Since the translation of each individual sequence function leaves the last operand empty, we append a ‘true’ (\top) as the final operand, since we don’t care what happens in the state after the sequence is complete.

(once ψ) := $\psi X \dots$

(hold ψ) := $\psi U \dots$

(hold-while $\psi \alpha \beta \dots \nu$) := $(\psi M \alpha) X (\psi M \beta) X \dots X (\psi M \nu) X \psi U \dots$ where the last $\psi U \dots$ allows for additional states satisfying ψ until the next modal is satisfied.

For example, a sequence such as the following, which signifies a throw attempt:

```
(then
  (once (agent_holds ?b))
  (hold (and (not (agent_holds ?b)) (in_motion ?b)))
  (once (not (in_motion ?b)))
)
```

Can be translated to LTL using $\psi := (\text{agent_holds } ?b)$, $\varphi := (\text{in_motion } ?b)$ as:

$\psi X (\neg \psi \wedge \varphi) U (\neg \varphi) X \top$

Here’s another example:

```
(then
  (once (agent_holds ?b)) ;  $\alpha$ 
  (hold-while
    (and (not (agent_holds ?b)) (in_motion ?b)) ;  $\beta$ 
    (touch ?b ?r) ;  $\gamma$ 
  )
  (once (and (in ?h ?b) (not (in_motion ?b)))) ;  $\delta$ 
)
```

If we translate each predicate to the letter appearing at the end of the line, this translates to:

$\alpha X (\beta M \gamma) X \beta U \delta X \top$

¹These are the ones we’ve used so far in the interactive experiment dataset, even if we previously defined other ones, too.